



Εθνικό Μετσόβιο Πολυτεχνείο
Σχολή Ηλεκτρολόγων Μηχανικών
και Μηχανικών Υπολογιστών
Τομέας Τεχνολογίας Πληροφορικής και Υπολογιστών

**Σχεδιασμός και Υλοποίηση Γεννήτριας Ρεαλιστικών
Χωροχρονικών Δεδομένων Μεγάλου Όγκου για Αποτίμηση
Υπηρεσιών Κοινωνικής Δικτύωσης**

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

ΘΑΛΕΙΑ-ΔΗΜΗΤΡΑ ΔΟΥΔΑΛΗ

Επιβλέπων : Νεκτάριος Κοζύρης
Καθηγητής Ε.Μ.Π.

Αθήνα, Ιούλιος 2015



Εθνικό Μετσόβιο Πολυτεχνείο
Σχολή Ηλεκτρολόγων Μηχανικών
και Μηχανικών Υπολογιστών
Τομέας Τεχνολογίας Πληροφορικής και Υπολογιστών

**Σχεδιασμός και Υλοποίηση Γεννήτριας Ρεαλιστικών
Χωροχρονικών Δεδομένων Μεγάλου Όγκου για Αποτίμηση
Υπηρεσιών Κοινωνικής Δικτύωσης**

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

ΘΑΛΕΙΑ-ΔΗΜΗΤΡΑ ΔΟΥΔΑΛΗ

Επιβλέπων : Νεκτάριος Κοζύρης
Καθηγητής Ε.Μ.Π.

Εγκρίθηκε από την τριμελή εξεταστική επιτροπή την 20η Ιουλίου 2015.

.....
Νεκτάριος Κοζύρης
Καθηγητής Ε.Μ.Π.

.....
Νικόλαος Παπασπύρου
Αν. Καθηγητής Ε.Μ.Π.

.....
Δημήτριος Τσουμάκος
Επίκουρος Καθηγητής
Ιόνιο Πανεπιστήμιο

Αθήνα, Ιούλιος 2015

.....
Θάλεια-Δήμητρα Δούδαλη

Διπλωματούχος Ηλεκτρολόγος Μηχανικός και Μηχανικός Υπολογιστών Ε.Μ.Π.

Copyright © Θάλεια-Δήμητρα Δούδαλη, 2015.

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι αντιπροσωπεύουν τις επίσημες θέσεις του Εθνικού Μετσόβιου Πολυτεχνείου.

Περίληψη

Η ποσότητα δεδομένων κοινωνικής δικτύωσης, η οποία παράγεται καθημερινά, αυξάνεται με ραγδαίους ρυθμούς. Η αποθήκευση και ο διαμοιρασμός του τεράστιου αυτού όγκου δεδομένων δε μπορεί πλέον να πραγματοποιηθεί με παραδοσιακές τεχνικές. Σαν αποτέλεσμα, οι σύγχρονες υπηρεσίες κοινωνικής δικτύωσης χρησιμοποιούν καταναμημένα συστήματα διαχείρισης δεδομένων, τα οποία τους παρέχουν επαρκή χώρο αποθήκευσης πληθώρας δεδομένων αλλά και μεθόδους ταχύτατης επεξεργασίας τους. Η αξιολόγηση των υπηρεσιών αυτών μπορεί να γίνει μέσω της μελέτης της επίδοσής τους κατά τη διάρκεια χρήσης τους. Η πλήρης κατανόηση και αποτίμηση, όμως, της υποδομής και των τεχνικών αποθήκευσης και επεξεργασίας δεδομένων, τις οποίες ακολουθούν, δε μπορούν να πραγματοποιηθούν λόγω της αδυναμίας πρόσβασης στον όγκο δεδομένων τον οποίο διαχειρίζονται. Αυτό οφείλεται στο γεγονός ότι πρόκειται για ιδιωτικά επιχειρησιακά δεδομένα που αφορούν πραγματικούς χρήστες και δεν μπορούν να εκμεταλευτούν ερευνητικά. Για το λόγο αυτό, σκοπός της παρούσας διπλωματικής εργασίας είναι η δημιουργία μίας γεννήτριας ρεαλιστικών χωροχρονικών δεδομένων μεγάλου όγκου, τα οποία θα προσομοιάζουν πραγματικά δεδομένα υπηρεσιών κοινωνικής δικτύωσης.

Πιο συγκεκριμένα, η γεννήτρια διαθέτει ως πηγή δεδομένων πραγματικά σημεία ενδιαφέροντος και κριτικές για τα σημεία αυτά από γνωστή υπηρεσία κοινωνικής δικτύωσης. Στη συνέχεια, δημιουργεί ημερήσιες ρεαλιστικές τροχιές χρηστών στο χάρτη χρησιμοποιώντας την υπηρεσία εύρεσης διαδρομών της Google. Για κάθε τροχιά αποθηκεύει τα δορυφορικά στίγματα των διαδρομών του χρήστη και τις επισκέψεις του στα σημεία ενδιαφέροντος, οι οποίες συνοδεύονται από βαθμολογία και κριτική του σημείου αυτού. Οι ημερήσιες αυτές τροχιές είναι διαθέσιμες σε μορφή στατικού χάρτη, όπως αυτός δημιουργείται από την αντίστοιχη υπηρεσία της Google, καθώς επίσης και με τη μορφή raw data. Η γεννήτρια λαμβάνει διάφορες παράμετρους εισόδου, οι οποίες διαφοροποιούν το συνολικά παραγόμενο όγκο και τη μορφή των δεδομένων. Ενδεικτικά, αυτές είναι το πλήθος χρηστών που θα δημιουργηθούν, το χρονικό διάστημα δημιουργίας ημερήσιων τροχιών καθώς και το πλήθος και η διάρκεια των ημερήσιων επισκέψεων.

Η ημερήσια εκτέλεση της γεννήτριας για ένα σημαντικό χρονικό διάστημα οδήγησε στη δημιουργία ενός συνόλου χωροχρονικών δεδομένων και δεδομένων κειμένου μεγάλου όγκου. Πιο συγκεκριμένα, δημιουργήθηκαν 9464 χρήστες, 1586537 επισκέψεις χρηστών και 38800019 δορυφορικά στίγματα, τα οποία αντιστοιχούν σε συνολικά δεδομένα μεγέθους 3 GB. Το σύνολο των δεδομένων αυτών αποθηκεύτηκε σε μία καταναμημένη βάση δεδομένων, στα πρότυπα αυτών που χρησιμοποιούν γνωστές πλατφόρμες κοινωνικής δικτύωσης. Στη συνέχεια, υλοποιήθηκαν επερωτήσεις στα διαθέσιμα δεδομένα, οι οποίες είναι αντιπροσωπευτικές πραγματικών επερωτήσεων σε αντίστοιχες πλατφόρμες. Τέλος, πραγματοποιήθηκε η εκτέλεση των επερωτήσεων αυτών στα δεδομένα αυτά για μεταβλητό πλήθος κόμβων του καταναμημένου συστήματος αποθήκευσης και μεταβλητό πλήθος ταυτόχρονων επερωτήσεων. Με τον τρόπο αυτό αξιολογήθηκε η κλιμακωσιμότητα του συστήματος αυτού, η οποία μας οδηγεί και σε μία ενδεικτική αποτίμηση των υπηρεσιών κοινωνικής δικτύωσης, οι οποίες χρησιμοποιούν αντίστοιχα καταναμημένα συστήματα αποθήκευσης και επεξεργασίας δεδομένων μεγάλου όγκου.

Λέξεις κλειδιά

Γεννήτρια, Χωροχρονικά δεδομένα, Δεδομένα μεγάλου όγκου, Σημεία ενδιαφέροντος, Καταναμημένα συστήματα, Υπηρεσίες Google, HBase, κλιμακωσιμότητα

Abstract

Nowadays, in the era of Big Data, the amount of social media data, that is being produced daily, increases significantly. The storage and analysis of such data cannot be achieved any more with traditional means and methods. Consequently, social networking services resort in using distributed systems and techniques, in order to store and manage effectively the huge amount of the data they own. Usually, the evaluation of such services results through the ease of use and satisfactory performance. However, deep understanding of how data is stored and managed cannot be reached, due to the lack of access to these data which is imposed by privacy restrictions. In this way, it is not possible to evaluate properly such social networking services. Therefore, goal of the current diploma thesis is to design and implement a generator of realistic spatio-temporal and textual data, that will be similar to real social media data.

The generator uses as source data, real points of interest and reviews for these points, extracted by a well-known travel service. Then, it creates realistic daily routes per user on the map, using the Google Directions API. These daily routes are available in the form of static maps, using the Google Static Maps API. Each daily route includes check-ins at points of interest, together with rating and review of the point, and gps traces indicating the route. Also, the generator functions with various input parameters, that differentiate the amount and structure of data produced. For example, such parameters can be the number of users created, the time period in which daily routes will be produced per user and the number and duration of daily check-ins.

The generator was executed daily for a significant amount of time in order to create a Big Data dataset of spatio-temporal and textual data. More specifically, the generator created 9464 users, 1586537 check-ins and 38800019 GPS traces, which sum up to 3 GB data. The dataset was stored in a distributed database system using a specific data storage model. Moreover, we implemented certain queries for these data, that are representative of queries imposed by the users of real social networking services. Finally, we created a workload of queries and executed them for different number of concurrent queries and different number of nodes of the distributed database system. In this way, we were able to perform a scalability testing to the system and evaluate the performance of distributed means of storage and processing of social media data used by many social networking services.

Keywords

Generator, Spatio-temporal data, Textual data, Big Data, Points of Interest, Daily routes, Google Directions API, Google Static Maps API, HBase, Scalability Testing, Distributed Systems

Ευχαριστίες

Η παρούσα διπλωματική εργασία εκπονήθηκε στο Εργαστήριο Υπολογιστικών Συστημάτων της Σχολής Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου, υπό την επίβλεψη του καθηγητή Νεκτάριου Κοζύρη.

Αρχικά, θα ήθελα να ευχαριστήσω τον καθηγητή μου Νεκτάριο Κοζύρη για την ευκαιρία που μου έδωσε να ασχοληθώ με ένα σύγχρονο και ενδιαφέρον θέμα, καθώς και για τη στήριξή του κατά την εκπόνηση της διπλωματικής εργασίας αλλά και τα επόμενα βήματά μου. Επίσης, θέλω να ευχαριστήσω το μεταδιδακτορικό ερευνητή Γιάννη Κωνσταντίνου για τη συνεχή καθοδήγηση και βοήθεια κατά τη διάρκεια της διπλωματικής εργασίας. Ακόμη, θέλω να ευχαριστήσω το διδακτορικό φοιτητή Γιάννη Γιαννακόπουλο για την αμέριστη βοήθεια σε πρακτικά ζητήματα σε όλη την πορεία εκπόνησης διπλωματικής εργασίας. Επιπρόσθετα, θέλω να ευχαριστήσω και τον αναπληρωτή καθηγητή Νικόλαο Παπασπύρου για τη βοήθειά του στην πορεία των σπουδών μου.

Επιπλέον, θέλω να ευχαριστήσω τους φίλους και συμφοιτητές μου για τις ευχάριστες στιγμές που περάσαμε τα 5 αυτά χρόνια. Τέλος, οφείλω να ευχαριστήσω την οικογένεια μου, τη μητέρα μου και τους αδερφούς μου Γιάννη και Στέλιο για τη συνεχή υποστήριξη και συμπαράσταση σε όλο το διάστημα των προπτυχιακών μου σπουδών.

Θάλεια-Δήμητρα Δούδαλη,

Αθήνα, 20η Ιουλίου 2015

Contents

1. Εισαγωγή	17
1.1 Κίνητρο	17
1.2 Κύρια σημεία της εργασίας	18
1.3 Οργάνωση κειμένου	18
1.4 Συνοπτική περιγραφή της γεννήτριας	18
1.5 Σύνοψη πειραματικής αξιολόγησης	19
1. Introduction	23
1.1 Motivation	23
1.2 Thesis contribution	23
1.3 Chapter outline	24
1.4 Brief generator description	24
2. Background	27
2.1 Databases	27
2.1.1 Relational databases	27
2.1.2 ACID properties	27
2.1.3 OLTP vs OLAP	28
2.2 PostgreSQL	28
2.2.1 PostGIS	29
2.2.2 Indexes	29
2.3 HDFS	30
2.4 HBase	30
2.4.1 Data model	30
2.4.2 Architecture components	31
2.4.3 Coprocessors	32
2.5 Google Maps Directions API	32
2.5.1 JSON response file	32
2.5.2 Encoded Polyline Algorithm Format	36
2.6 Google Maps Static Maps API	36
2.7 Mathematical Background	37
2.7.1 Random Variables	37
2.7.2 Uniform Distribution	37
2.7.3 Normal Distribution	38
2.7.4 Bernoulli Distribution	40
3. Generator	41
3.1 Storage of source data	41
3.2 Generator attributes	42
3.2.1 User	42
3.2.2 Check-in	42
3.2.3 GPS trace	43

3.2.4	Static route map - MapURL	43
3.2.5	Point of interest - Poi	43
3.2.6	Review	43
3.3	Input parameters	43
3.4	Generator's configuration parameters	44
3.4.1	Home	44
3.4.2	Trip	45
3.4.3	Check-in POI	45
3.4.4	Check-in review	46
3.4.5	Path between check-ins	46
3.4.6	Timestamp	46
3.4.7	Static map	47
4.	Dataset	49
4.1	Generator input parameters	49
4.2	Generator deployment setup	49
4.3	Dataset storage model	50
4.3.1	Friends Table	50
4.3.2	Check-ins Table	51
4.3.3	GPS traces Table	51
4.4	Dataset storage architecture	52
5.	Queries	53
5.1	Query implementation	53
5.1.1	Most visited POI query	53
5.1.2	News Feed query	54
5.1.3	Correlated Most Visited POI query	54
5.2	Queries execution	55
5.2.1	Workload	55
5.2.2	Workload generation setup	55
5.2.3	Scalability testing	55
6.	Conclusion	59
6.1	Summary	59
6.2	Related Work	59
6.3	Future Work	60
	Bibliography	61

List of Figures

1.1	Μετρήσεις κλιμακωσιμότητας	20
2.1	A table in HBase consisting of two column families	31
2.2	Discrete uniform probability mass function	38
2.3	Discrete uniform cumulative distribution function	38
2.4	Probability density function	39
2.5	Cumulative distribution function	39
2.6	3-sigma empirical rule for the normal distribution	39
3.1	Database schema for source data	41
3.2	Class diagram of generator's attributes	42
3.3	Example of static map image	47
4.1	Generator execution architecture	50
4.2	HBase cluster architecture	52
5.1	Query execution architecture	55
5.2	Scalability testing - Latency	57
5.3	Scalability testing - Throughput	58

List of Tables

4.1	Dataset produced by the generator	50
4.2	Example of HBase Table 'friends'	51
4.3	Example of HBase Table 'check-ins'	51
4.4	Example of HBase Table 'gps-traces'	51
5.1	Scalability testing - Latency	56
5.2	Scalability testing - Throughput	57

Chapter 1

Εισαγωγή

1.1 Κίνητρο

Στη σύγχρονη εποχή, τα δεδομένα κοινωνικής δικτύωσης, τα οποία παράγονται και διαχειρίζονται καθημερινά από τις αντίστοιχες υπηρεσίες, αυξάνονται με ραγδαίους ρυθμούς. Επιπλέον, τα δεδομένα αυτά περιέχουν πληροφορίες κείμενου, φωτογραφίας, βίντεο, χρόνου και τοποθεσίας. Ενδεικτικά, η διαδικτυακή υπηρεσία κοινωνικής δικτύωσης Twitter, σύμφωνα με τα επίσημα στατιστικά [1], διαθέτει 302 εκατομμύρια ενεργούς χρήστες κάθε μήνα και διαχειρίζεται κάθε μέρα 500 εκατομμύρια μηνύματα (tweets) τα οποία δημοσιεύονται από τους χρήστες. Τα tweets περιλαμβάνουν στοιχεία κειμένου και χρόνου (η ώρα δημοσίευσης) και μπορεί να περιέχουν επίσης εικόνα, βίντεο και πληροφορίες τοποθεσίας. Τα παραδοσιακά μέσα αποθήκευσης και τεχνικές επεξεργασίας δεν επαρκούν πλέον για τον τεράστιο όγκο των πολύμορφων κοινωνικών δεδομένων. Για το λόγο αυτό, οι σύγχρονες υπηρεσίες κοινωνικής δικτύωσης καταφεύγουν σε καταναμεμημένα συστήματα και εργαλεία διαχείρισης δεδομένων. Για παράδειγμα, η υπηρεσία Twitter αρχικά αποθήκευε τα tweets των χρηστών σε παραδοσιακές βάσεις δεδομένων [5]. Όμως, με το πέρασμα του καιρού και την τεράστια αύξηση των tweets, υπήρχαν προβλήματα στην ανάγνωση και εγγραφή δεδομένων στις βάσεις αυτές. Έτσι, η υπηρεσία δημιούργησε και χρησιμοποίησε καταναμεμημένα εργαλεία αποθήκευσης για τα δεδομένα τα οποία έπρεπε να διαχειριστεί [7]. Με τα εργαλεία αυτά παρατηρήθηκε σημαντική βελτίωση στην επίδοση της υπηρεσίας.

Ένας μεγάλος αριθμός από διαδικτυακές υπηρεσίες χρησιμοποιούν πλέον καταναμεμημένα συστήματα αποθήκευσης δεδομένων. Για παράδειγμα, μεγάλες διαδικτυακές υπηρεσίες, όπως το Twitter και το Yahoo!, χρησιμοποιούν το καταναμεμημένο σύστημα βάσης δεδομένων HBase, για τη διαχείριση ενός τμήματος των δεδομένων τους [7]. Αντίστοιχα, πολλές διαδικτυακές υπηρεσίες κοινωνικής δικτύωσης, όπως το Facebook και το Twitter, χρησιμοποιούν το καταναμεμημένο σύστημα διαχείρισης δεδομένων Hadoop [8]. Η ελεύθερη γνώση των συστημάτων τα οποία χρησιμοποιούν οι υπηρεσίες αυτές, λόγω του γεγονότος ότι είναι συστήματα ανοιχτού κώδικα, μας δίνει μία πρώτη εικόνα της λειτουργίας τους. Όμως, τα δεδομένα τα οποία διαχειρίζονται δε γίνονται διαθέσιμα στο ευρύ κοινό, για λόγους διαφύλαξης των προσωπικών δεδομένων των χρηστών. Έτσι, δε μπορούμε να κατανοήσουμε πλήρως τη λειτουργία των υπηρεσιών αυτών και να αξιολογήσουμε την επιλογή των συγκεκριμένων καταναμεμημένων εργαλείων για τη διαχείριση των δεδομένων τους.

Με την παρούσα διπλωματική εργασία, ερχόμαστε να καλύψουμε αυτό το κενό στην πλήρη κατανόηση της διαχείρισης δεδομένων χρηστών κατά την λειτουργία των υπηρεσιών κοινωνικής δικτύωσης. Εφ'όσον δε μπορούμε να έχουμε πρόσβαση στα αντίστοιχα δεδομένα, προσπαθήσαμε να δημιουργήσουμε ρεαλιστικά δεδομένα κοινωνικής δικτύωσης, τα οποία να προσομοιάζουν σε μεγάλο βαθμό αντίστοιχα πραγματικά δεδομένα. Για το σκοπό αυτό, σχεδιάσαμε και υλοποιήσαμε μία γεννήτρια ρεαλιστικών χωροχρονικών δεδομένων, στα πρότυπα αυτών που συναντάμε σε μεγάλες υπηρεσίες κοινωνικής δικτύωσης, όπως το Facebook.

1.2 Κύρια σημεία της εργασίας

Οι βασικές συνεισφορές της διπλωματικής εργασίας είναι οι ακόλουθες:

1. Σχεδιασμός και υλοποίηση παραμετροποιήσιμης γεννήτριας χωροχρονικών δεδομένων ανοιχτού κώδικα ¹.
2. Δημιουργία συνόλου χωροχρονικών δεδομένων μεγάλου όγκου με χρήση της γεννήτριας.
3. Αποδοτική αποθήκευση και δεικτοδότηση του συνόλου δεδομένων σε κατανεμημένο σύστημα διαχείρισης δεδομένων.
4. Πειραματική αξιολόγηση της κλιμακωσιμότητας του συστήματος διαχείρισης δεδομένων, κάτω από διαφορετικό όγκο δεδομένων, ταυτόχρονα συνδεδεμένων χρηστών και αριθμών κόμβων.

1.3 Οργάνωση κειμένου

Στη συνέχεια του κεφαλαίου 1 ακολουθεί μία σύντομη περιγραφή του σχεδιασμού της γεννήτριας χωροχρονικών δεδομένων. Επίσης, αναφέρονται συνοπτικά τα αποτελέσματα της αποτίμησης των υπηρεσιών κοινωνικής δικτύωσης μέσω της υλοποίησης και εφαρμογής επερωτήσεων στο σύνολο δεδομένων μεγάλου όγκου, τα οποία δημιουργήθηκαν από την γεννήτρια και αποθηκεύτηκαν σε ένα σύστημα κατανεμημένης διαχείρισης δεδομένων.

Στο κεφάλαιο 2 παρουσιάζεται το θεωρητικό υπόβαθρο, το οποίο είναι σημαντικό να γνωρίζει κάποιος, ώστε να κατανοήσει έννοιες και ορολογίες που θα αναφέρονται στη συνέχεια. Ενδεικτικά, αναλύονται τα εργαλεία που χρησιμοποιήθηκαν στην υλοποίηση της γεννήτριας, όπως η βάση δεδομένων για τα πηγαία δεδομένα, οι υπηρεσίες της Google για τη δημιουργία και παρουσίαση των τροχιών και ορισμένες μαθηματικές έννοιες για τις παραμέτρους εισόδου. Επίσης, παρουσιάζονται τα κατανεμημένα εργαλεία αποθήκευσης και διαχείρισης δεδομένων, τα οποία χρησιμοποιήθηκαν για τα παραγόμενα από τη γεννήτρια δεδομένα.

Στο κεφάλαιο 3 παρατίθεται μία λεπτομερής περιγραφή της μεθοδολογίας σχεδιασμού της γεννήτριας χωροχρονικών δεδομένων. Αναλύεται ο τρόπος αποθήκευσης των πηγαίων δεδομένων, ορίζονται οι παράμετροι εισόδου της γεννήτριας και περιγράφονται με λεπτομέρεια οι οντότητες, οι αποφάσεις, οι παραδοχές και η πλήρης λειτουργία της γεννήτριας.

Στο κεφάλαιο 4 αναλύεται ο τρόπος δημιουργίας ενός συνόλου δεδομένων μεγάλου όγκου παραγόμενων από τη γεννήτρια, καθώς και το σχήμα δεδομένων για την αποθήκευσή του σε μία κατανεμημένη βάση δεδομένων.

Στο κεφάλαιο 5 παρατίθεται η μεθοδολογία υλοποίησης των επερωτήσεων στο σύνολο των διαθέσιμων δεδομένων τα οποία παράχθηκαν από τη γεννήτρια. Επίσης, παρουσιάζονται τα αποτελέσματα από την αξιολόγηση της κλιμακωσιμότητας του συστήματος της βάσης δεδομένων στην οποία αποθηκεύτηκαν τα δεδομένα μέσω της εκτέλεσης των επερωτήσεων αυτών στα δεδομένα.

Στο κεφάλαιο 6 ανακεφαλαιώνονται τα συμπεράσματα της παραπάνω αποτίμησης των εργαλείων αποθήκευσης και επεξεργασίας δεδομένων που χρησιμοποιούν οι υπηρεσίες κοινωνικής δικτύωσης. Επίσης, παρουσιάζονται σχετικές εργασίες και πιθανές μελλοντικές επεκτάσεις της τρέχουσας διπλωματικής.

1.4 Συνοπτική περιγραφή της γεννήτριας

Η θεμελιώδης λειτουργία της γεννήτριας είναι η παραγωγή ρεαλιστικών χωροχρονικών δεδομένων. Πιο συγκεκριμένα, η γεννήτρια δημιουργεί χρήστες υπηρεσιών κοινωνικής δικτύωσης, οι οποίοι κάθε μέρα επισκέπτονται διάφορα σημεία ενδιαφέροντος για τα οποία αφήνουν μία κριτική

¹ Ο κώδικας είναι διαθέσιμος στο <https://github.com/Thaleia-DimitraDoudali/thesis>

και βαθμολογία στην πλατφόρμα. Επομένως, η γεννήτρια παράγει ημερήσιες τροχιές χρηστών για ένα καθορισμένο χρονικό διάστημα.

Η γεννήτρια ρεαλιστικών χωροχρονικών δεδομένων χρησιμοποιεί πραγματικά σημεία ενδιαφέροντος ως πηγή δεδομένων. Τα σημεία αυτά προμηθεύτηκαν από την υπηρεσία ταξιδιών κοινωνικής δικτύωσης TripAdvisor. Τα πηγαία δεδομένα περιέχουν τις γεωγραφικές συντεταγμένες, την ονομασία και διεύθυνση των σημείων ενδιαφέροντος, καθώς και μία λίστα από κριτικές και βαθμολογίες για τα αντίστοιχα σημεία από πραγματικούς χρήστες της υπηρεσίας. Για την αποθήκευση των δεδομένων αυτών χρησιμοποιήθηκε η βάση δεδομένων PostgreSQL, έτσι ώστε να έχουμε πρόσβαση στις συναρτήσεις χειρισμού γεωγραφικών τύπων δεδομένων και υπολογισμού αποστάσεων μεταξύ γεωγραφικών σημείων, οι οποίες παρέχονται από την επέκταση PostGIS την οποία διαθέτει η PostgreSQL. Η γεννήτρια λαμβάνει ορισμένες παραμέτρους εισόδου, έτσι ώστε να παραμετροποιείται η λειτουργία της και να τροποποιείται το πλήθος και εν μέρει η δομή των παραγόμενων δεδομένων. Ενδεικτικά, οι παράμετροι αφορούν το πλήθος των χρηστών που θα παραχθούν, το χρονικό διάστημα για το οποίο θα δημιουργηθούν ημερήσιες τροχιές για όλους τους χρήστες και ορισμένες μετρικές που επηρεάζουν τις ημερήσιες επισκέψεις των χρηστών στα σημεία ενδιαφέροντος.

Η γεννήτρια κατά τη δημιουργία των ημερήσιων τροχιών λαμβάνει αποφάσεις οι οποίες βασίζονται σε παράγοντες τυχαιότητας. Πιο συγκεκριμένα, τα σημεία ενδιαφέροντος τα οποία επισκέπτεται κάθε χρήστης καθορίζονται με τυχαίο τρόπο από τη γεννήτρια, με μόνο περιορισμό να είναι σε κοντινή απόσταση μεταξύ τους και από την κατοικία του χρήστη. Η κατοικία κάθε χρήστη ορίζεται και αυτή με τυχαίο τρόπο και η σχετική απόσταση καθορίζεται από παράμετρο εισόδου. Οι χρήστες έχουν τη δυνατότητα να πραγματοποιούν ταξίδια, έτσι ώστε να επισκέπτονται και σημεία ενδιαφέροντος σε πιο μακρινή απόσταση. Η τοποθεσία και διάρκεια του ταξιδιού καθορίζονται και αυτές με τυχαίο τρόπο. Η γεννήτρια επιτρέπει σε ένα χρήστη να επισκέπτεται κάποιο σημείο μόνο μία φορά την ημέρα. Σε κάθε σημείο στο οποίο ο χρήστης κάνει check-in, η γεννήτρια αναθέτει μία τυχαία κριτική και βαθμολογία για το σημείο αυτό από το συγκεκριμένο χρήστη. Η κριτική αυτή και η βαθμολογία του σημείου αυτού είναι πραγματική, όπως την έγραψε κάποιος χρήστης της υπηρεσίας TripAdvisor.

Σχετικά με τη διαδρομή από ένα σημείο ενδιαφέροντος στο επόμενο, ο χρήστης περπατά για να μεταβεί στον προορισμό του, γι' αυτό και η απόσταση μεταξύ των σημείων έχει νόημα να είναι κοντινή και καθορίζεται από αντίστοιχη παράμετρο εισόδου. Η γεννήτρια χρησιμοποιεί την υπηρεσία εύρεσης διαδρομών της Google (Google Directions API) [2]. Στέλνει μία αίτηση στην υπηρεσία, η οποία τελικά επιστρέφει ένα αρχείο απάντησης με αναλυτικές οδηγίες και πληροφορίες για τη διαδρομή. Από το αρχείο αυτό η γεννήτρια αξιοποιεί το σύνολο των σημείων της διαδρομής, τα οποία επιστρέφει η υπηρεσία, ως στίγματα δορυφόρου που υποδεικνύουν τη διαδρομή. Η γεννήτρια λαμβάνει, επίσης, από το αρχείο απάντησης και τη διάρκεια της διαδρομής από το σημείο αρχής στο σημείο προορισμού. Με την πληροφορία αυτή, μπορεί να υπολογίσει και τη χρονική στιγμή την οποία ο χρήστης θα κάνει check-in στο επόμενο σημείο. Το πλήθος των σημείων ενδιαφέροντος τα οποία θα επισκεφτεί ο χρήστης κατά τη διάρκεια της ημέρας καθορίζονται με τυχαίο τρόπο, όπως και η διάρκεια της επίσκεψης σε κάθε σημείο. Τέλος, η γεννήτρια χρησιμοποιεί την υπηρεσία στατικού χάρτη της Google (Google Static Maps API) [4], προκειμένου να απεικονίσει στο χάρτη τις ημερήσιες διαδρομές των χρηστών.

Συμπερασματικά, τα δεδομένα τα οποία παράγει η γεννήτρια προσομοιάζουν σε μεγάλο βαθμό πραγματικά δεδομένα κοινωνικής δικτύωσης. Αυτό οφείλεται στη χρήση πραγματικών σημείων ενδιαφέροντος ως πηγαία δεδομένα, πραγματικών κριτικών που αναφέρονται στα δεδομένα αυτά, ρεαλιστικών τροχιών και διαρκειών όπως αυτές παράγονται από το Google Directions API, καθώς και στην αξιοποίηση παραγόντων τυχαιότητας στις αποφάσεις τις οποίες λαμβάνει η γεννήτρια.

1.5 Σύνοψη πειραματικής αξιολόγησης

Η γεννήτρια εκτελέστηκε για ένα σημαντικό χρονικό διάστημα με σκοπό τη δημιουργία ενός συνόλου χωροχρονικών δεδομένων και δεδομένων κειμένου μεγάλου όγκου. Πιο συγκεκριμένα, δημιουργήθηκαν 9464 χρήστες, 1586537 επισκέψεις χρηστών και 38800019 δορυφορικά στίγματα,

τα οποία αντιστοιχούν σε συνολικά δεδομένα μεγέθους 3 GB. Στα δεδομένα αυτά, προσθέσαμε και έναν γράφο φίλων της υπηρεσίας κοινωνικής δικτύωσης Twitter τον οποίο είχαμε στη διάθεσή μας. Τα δεδομένα αυτά αποθηκεύτηκαν με ένα συγκεκριμένο σχήμα στην κατανεμημένη βάση.

Στη συνέχεια, υλοποιήθηκαν συγκεκριμένες επερωτήσεις στα δεδομένα αυτά με εργαλεία που αξιοποιούν αποδοτικά όλους τους κόμβους του κατανεμημένου συστήματος. Πιο συγκεκριμένα, οι επερωτήσεις αυτές είναι οι ακόλουθες:

1. Για ένα τυχαίο χρήστη της υπηρεσίας, βρες το σημείο ενδιαφέροντος το οποίο επισκεφτηκε τις περισσότερες φορές.
2. Για ένα τυχαίο χρήστη της υπηρεσίας, εμφάνισε κατα χρονολογική σειρά τις επισκέψεις όλων των φίλων του για μία συγκεκριμένη ημέρα.
3. Για ένα τυχαίο χρήστη της υπηρεσίας, βρες ποιοι φίλοι του και πόσες φορές έχουν επισκεφτεί το σημείο το οποίο ο ίδιος έχει επισκεφτεί τις περισσότερες φορές.

Εκτελέσαμε τις επερωτήσεις αυτές με ταυτόχρονο ρυθμό για μεταβλητό συνολικό πλήθος επερωτήσεων. Επίσης, μεταβάλλαμε τον αριθμό κόμβων της κατανεμημένης βάσης δεδομένων και με τον τρόπο αυτό αξιολογήσαμε την κλιμακωσιμότητα του συστήματος βάσης. Τα αποτελέσματα της εκτέλεσης αυτής συνοψίζονται στα ακόλουθα διαγράμματα:

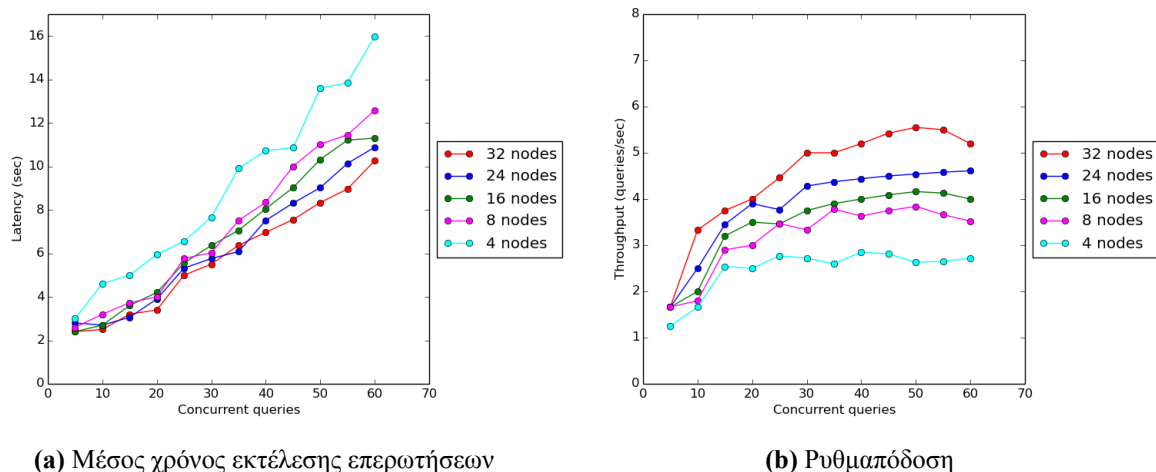


Figure 1.1: Μετρήσεις κλιμακωσιμότητας

Από τα παραπάνω διαγράμματα παρατηρούμε ότι ο μέσος χρόνος εκτέλεσης των επερωτήσεων είναι μεγαλύτερος όσο αυξάνεται το πλήθος των ταυτόχρονων επερωτήσεων. Αυτό είναι λογικό από τη στιγμή που οι εξυπηρετητές πρέπει να διαχειριστούν και να διεκπεραιώσουν περισσότερα αιτήματα, οπότε κάποια αιτήματα θα περιμένουν στις ουρές των εξυπηρετητών όσο εκείνοι διεκπεραιώνουν προηγούμενα αιτήματα. Για το λόγο αυτό, παρατηρούμε ότι για διαφορετικό αριθμό κόμβων ο μέσος χρόνος εκτέλεσης μικρού πλήθους ταυτόχρονων επερωτήσεων είναι παραπλήσιος, καθώς αυτά απαντούνται άμεσα χωρίς να περιμένουν σε ουρές αναμονής. Επίσης, ο μέσος χρόνος εκτέλεσης αυξάνεται όσο μειώνεται το πλήθος των κόμβων. Και αυτό είναι αναμενόμενο, καθώς μικρότερο πλήθος κόμβων δε μπορεί να εξυπηρετήσει τα ίδια και περισσότερα αιτήματα με τον ίδιο ρυθμό.

Όσον αφορά τη ρυθμαπόδοση του συστήματος, παρατηρούμε ότι αυτή μειώνεται όσο αυξάνεται το πλήθος των κόμβων. Αυτό δικαιολογείται από το γεγονός ότι λιγότεροι κόμβοι μπορούν να εξυπηρετήσουν τα αιτήματα με πιο αργό ρυθμό καθώς δέχονται περισσότερα ταυτόχρονα αιτήματα από ότι αν υπήρχαν περισσότεροι εξυπηρετητές. Παρ' όλα αυτά βλέπουμε ότι υπάρχει μία μέγιστη τιμή στην οποία μπορεί να φτάσει η ρυθμαπόδοση για κάθε πλήθος κόμβων. Αυτό συναντάται

σε όλα τα αντίστοιχα πραγματικά συστήματα και δείχνει ότι παρ'όλο που υπάρχουν περισσότεροι εξυπηρετητές, υπάρχει ένα άνω όριο στον αριθμό των επερωτήσεων που μπορούν να διεκπεραιώσουν στη μονάδα του χρόνου. Οι μικρές διακυμάνσεις στη μέγιστη αυτή τιμή οφείλονται στο γεγονός ότι οι επερωτήσεις αφορούν τυχαίο μέρος των δεδομένων. Επίσης, μπορεί να οφείλονται στο γεγονός ότι οι εξυπηρετητές χρησιμοποιούν και κρυφή μνήμη για τα δεδομένα που διαβάζουν από το δίσκο, με την ανάγνωση από την κρυφή μνήμη να είναι πολύ ταχύτερη. Επιπρόσθετα, το σύστημα κατανεμημένης μνήμης φιλοξενείται σε εικονικές μηχανές ενός υπολογιστικού νέφους, με αποτέλεσμα η λειτουργία του να επηρεάζεται από τους υπόλοιπους ταυτόχρονους χρήστες του νέφους.

Εν κατακλείδι, συμπεραίνουμε ότι το σύστημα κατανεμημένης μνήμης είναι κλιμακώσιμο για το συγκεκριμένο σχήμα αποθήκευσης των παραγόμενων από τη γεννήτρια δεδομένων μεγάλου όγκου, καθώς παρουσιάζει τη χαρακτηριστική συμπεριφορά ενός κλιμακώσιμου συστήματος.

Chapter 1

Introduction

1.1 Motivation

Nowadays, the amount of social networking data that are being produced and consumed daily is huge and it is constantly increasing. Moreover, these data can be in the form of text, photographs, video, time and location. For example, the online social networking service Twitter has 302 million active users per month and manages 500 million tweets per day, according to official statistics [1]. The tweets include textual data and a timestamp (the time of publication) and may also contain image, video, and information about location. Traditional data storage and processing techniques are no longer sufficient for the huge amount of polymorphic social media data. Consequently, modern social networking services employ distributed systems and data management tools. For example, in the beginning Twitter [5] used to store tweets in traditional databases. However, as time passed and the amount of tweets increased significantly, Twitter had several issues in reading and writing data in these databases. Thus, Twitter created and used distributed storage tools in order to manage its data. These tools gave notable boost in the performance of the service.

A large number of online services now use distributed data storage systems. For example, large online social media platforms such as Twitter and Yahoo!, use the distributed database system HBase, so as to manage their data [7]. Similarly, many online social networking services like Facebook and Twitter, use the distributed data storage and management system Hadoop [8]. Even though we can retrieve a list of social media platforms that use certain distributed data management tools, we cannot have access to the data that these platforms use, due to user's privacy restrictions. In this way, it is not possible to evaluate properly such social networking services, in terms of scalability and performance.

In the current diploma thesis, we want to bridge the gap considering the full understanding of the function of social networking services. Since we can not have access to the relevant data, we tried to create realistic social networking data, which are similar to real such data. Hence, we designed and implemented a generator of realistic spatiotemporal and textual data, similar to those found in well-known social networking services such as Facebook. The current diploma thesis is inspired and builds upon MoDisSENSE, which is a distributed platform that provides personalized search for points of interest and trending events based on the user's social graph by combining spatio-textual user generated data.[19]

1.2 Thesis contribution

The main contributions of this work are the following:

1. Design and implementation of an open source¹ parameterized generator of spatio-temporal and textual social media data.
2. Creation of a large dataset of such complex realistic social media data using the generator.

¹ Code available at <https://github.com/Thaleia-DimitraDoudali/thesis>

3. Creation of the appropriate data model and utilization of HBase, a distributed data management NoSQL database, for the storage and management of these data.
4. Scalability testing of the HBase cluster for the specific data storage model of the dataset produced by the generator.

1.3 Chapter outline

In chapter 2 we present the theoretical background, which is important for someone to know, in order to understand concepts and terminologies used later on. More specifically, we analyze the tools used in the implementation of the generator, such as the database for the source data, Google's services for the creation and presentation of user's routes and a number of mathematical concepts for the input parameters. Also, we include the distributed storage and data management tools, which were used in order to store the dataset created by the generator.

In chapter 3 we describe in detail the design of the generator. We analyze the storage schema of the source data, we define the input parameters and present thoroughly the attributes, the decision and complete function of the generator.

In chapter 4 we analyze the creation of a Big Data dataset using the generator and the data storage model according to which the dataset is stored into an HBase cluster.

In chapter 5 we present the implementation of several queries over the available dataset, as well as the creation of a workload consisting of these queries. Also, we perform the scalability testing of the HBase cluster and present the corresponding results.

In chapter 6 we sum up the conclusions from the above evaluation of the distributed storage and processing tools of the social media data. Finally, we present related work and possible future work.

1.4 Brief generator description

The main function of the generator is to create realistic spatiotemporal and textual social media data. More specifically, the generator creates users of online social networking services, who visit many places during the day and leave a review and rating to the corresponding points of interest. Therefore, the generator produces user's daily routes during a specific period of time.

The generator uses real points of interest as source data. These points were crawled from TripAdvisor, a well-known online travel service. The source data contain the geographical coordinates, the title and address of each point of interest, as well as a list of ratings and reviews for each point made by real TripAdvisor users. In order to store these points, we used the PostgreSQL database, so as to have access to functions of geographic data types, such as the calculation of distance, given by the extension PostGIS of PostgreSQL. The generator is parameterized by certain input variables, so as to create data that can vary in size and structure. For example, there are input parameters that define the number of users created, the period of time for which the generator will create daily routes and others that influence the daily check-ins to the points of interest.

The generator takes certain decisions, that refer to the daily routes created, using random factors. More specifically, the points of interest that a user visits are defined in a random way. The only restriction to this decision, is that the points must be near the user's home. The home of the user is defined in a random way and the distance from there is determined by an input parameter. The generator gives the ability to travel to its users, so that they can check-in into places that are further away from their home. The location and duration of their trips are also decided in a random way. The generator allows only one visit per day to each point of interest per user. Every user that visits a point of interest leaves a random rating and review for that point. This rating and review is real as it is written by real users of TripAdvisor.

As far as the transition from one point of interest to the next one is concerned, the user walks in order to get to his destination. This is the reason why the check-ins should occur in points of interest

that are in walking distance between them. The distance between the visits is defined through an input parameter. The generator issues a request to Google's Directions API [2] and gets a response file that contains analytical directions and information of the route from the origin to the destination. The response file includes all the intermediate points on the map in an encoded representation. The generator decodes and transforms these points to GPS traces, that indicate the route. The response file also contains the duration of the route from the starting point to the destination. Using this information, the generator can calculate the time when the user will check-in to the next place. Each check-in and GPS trace is timestamped. The number of points that a user will visit during the day and the duration of the visits is determined in a random way influenced by certain input parameters. Finally, the generator uses the Google Static Maps API [4], in order to illustrate each user's daily routes into a static map.

As a result, the data created by the generator simulate real social media data, due to the usage of real points of interest as source data, real reviews and rating to these points, real routes as given by the Google Directions API, and the utilization of random factors in the decisions that the generator makes, while creating the daily routes.

Chapter 2

Background

2.1 Databases

A database [9] is an organized collection of data, which can be managed in an easy way. A database management system is a computer software application that interacts with the user and other applications, in order to create, access, manage and execute queries at the data stored in the database. Databases are used in a wide range of services, such as bank, university and communication platforms.

2.1.1 Relational databases

A relational database is a digital database whose organization is based on the relational model of data [13]. According to this model, data are organized in tables with rows and columns. Each table represents a relation variable, each row is an implementation of this variable and each column is an attribute of the variable. Tables can be associated with one-to-one, one-to-many and many-to-many relationships. More specifically, a relation is a set of table rows that have common attributes. Moreover, each row can be uniquely identified by a certain attribute, called primary key. When a primary key is common between two tables, then it becomes a foreign key for the second table. Finally, a database index is a data structure that improves the speed of data retrieval operations on a database table at the cost of additional writes and storage space to maintain the index data structure.

2.1.2 ACID properties

ACID (Atomicity, Consistency, Isolation, Durability) [15] is a set of properties that guarantee that database transactions are processed reliably. In the context of databases, a single logical operation on the data is called a transaction.

Atomicity

Atomicity requires that each transaction be "all or nothing": if one part of the transaction fails, the entire transaction fails, and the database state is left unchanged. An atomic system must guarantee atomicity in each and every situation, including power failures, errors, and crashes. To the outside world, a committed transaction appears (by its effects on the database) to be indivisible ("atomic"), and an aborted transaction does not happen.

Consistency

The consistency property ensures that any transaction will bring the database from one valid state to another. Any data written to the database must be valid according to all defined rules, including constraints, cascades, triggers, and any combination thereof. This does not guarantee correctness of the transaction in all ways the application programmer might have wanted (that is the responsibility of application-level code) but merely that any programming errors cannot result in the violation of any defined rules.

Isolation

The isolation property ensures that the concurrent execution of transactions results in a system state that would be obtained if transactions were executed serially. Providing isolation is the main goal of concurrency control. Depending on concurrency control method (i.e. if it uses strict - as opposed to relaxed - serializability), the effects of an incomplete transaction might not even be visible to another transaction.

Durability

Durability means that once a transaction has been committed, it will remain so, even in the event of power loss, crashes, or errors. In a relational database, for instance, once a group of SQL statements execute, the results need to be stored permanently (even if the database crashes immediately thereafter). To defend against power loss, transactions (or their effects) must be recorded in a non-volatile memory.

2.1.3 OLTP vs OLAP

OLTP

Online transaction processing, or OLTP, is a class of information systems that facilitate and manage transaction-oriented applications, typically for data entry and retrieval transaction processing. OLTP is characterized by a large number of short on-line transactions (INSERT, UPDATE, DELETE). OLTP must possess ACID qualities to maintain data integrity and to ensure that transactions are correctly executed. The main emphasis for OLTP systems is put on very fast query processing, maintaining data integrity in multi-access environments and an effectiveness measured by number of transactions per second. In OLTP database there is detailed and current data, and the schema used to store transactional databases is the entity model.

OLAP

Online analytical processing, or OLAP, is characterized by relatively low volume of transactions. Queries are often very complex and involve aggregations. For OLAP systems a response time is an effectiveness measure. OLAP applications are widely used by Data Mining techniques. In OLAP database there is aggregated, historical data, stored in multi-dimensional schemas.

2.2 PostgreSQL

PostgreSQL [21] is an open source object-relational database management system. It runs on all major operating systems, including Linux, UNIX (AIX, BSD, HP-UX, SGI IRIX, Mac OS X, Solaris, Tru64), and Windows. It is fully ACID compliant, has full support for foreign keys, joins, views, triggers, and stored procedures (in multiple languages). It includes most SQL:2008 data types, including INTEGER, NUMERIC, BOOLEAN, CHAR, VARCHAR, DATE, INTERVAL, and TIMESTAMP. It also supports storage of binary large objects, including pictures, sounds, or video. Moreover, it has native programming interfaces for C/C++, Java, .Net, Perl, Python, Ruby, Tcl, ODBC. It supports international character sets, multibyte character encodings, Unicode, and it is locale-aware for sorting, case-sensitivity, and formatting. In addition, it supports unlimited database size, rows and indexes per table, 32TB table size, 1.6TB row size and 1GB field size. PostgreSQL manages database access permissions using the concept of roles. A role can be thought of as either a database user, or a group of database users, depending on how the role is set up. Roles can own database objects (for example, tables) and can assign privileges on those objects to other roles to control who has access to which objects. Furthermore, it is possible to grant membership in a role to another role, thus allowing the

member role use of privileges assigned to the role it is a member of. Finally, PostgreSQL comes with several extensions that add extra capabilities in its function and usage.

2.2.1 PostGIS

PostGIS [20] is a spatial database extender for PostgreSQL object-relational database. It adds support for geographic objects allowing location queries to be run in SQL. It defines new data types, functions, operators and indexes especially for geographic objects.

Geography type

The geography type provides native support for spatial features represented on geographic coordinates. A specific point on a map can be identified by its geographic coordinates, stated in the form of (latitude, longitude). Geographic coordinates are spherical coordinates expressed in angular units (degrees). The basis for the PostGIS geographic type is a sphere. The shortest path between two points on the sphere is a great circle arc. That means that calculations on geographies (areas, distances, lengths, intersections, etc) must be calculated on the sphere, using more complicated mathematics. For more accurate measurements, the calculations must take the actual spheroidal shape of the world into account. There are several functions and operators that take as input or return as output a geography data type object. One very useful one, is the function `ST_DWithin`. This function returns true if the geographies are within the specified distance of one another. Units are in meters and measurement is defaulted to measure around spheroid.

2.2.2 Indexes

A database index is a data structure that improves the speed of data retrieval operations on a database table at the cost of additional writes and storage space to maintain the index data structure. Indexes are used to quickly locate data without having to search every row in a database table every time a database table is accessed. Indexes can be created using one or more columns of a database table, providing the basis for both rapid random lookups and efficient access of ordered records.

B-tree

PostgreSQL uses by default a B-tree index on a table column. B-tree [10] is a tree data structure that keeps data sorted and allows searches, sequential access, insertions, and deletions in logarithmic time. The B-tree is a generalization of a binary search tree in that a node can have more than two children. More specifically, it is a balanced tree whose nodes contain pointers to a table's records and a number of keys. The keys act as separation values which divide its subtrees. For example, if an internal node has 3 child nodes (or subtrees) then it must have 2 keys: `a1` and `a2`. All values in the leftmost subtree will be less than `a1`, all values in the middle subtree will be between `a1` and `a2`, and all values in the rightmost subtree will be greater than `a2`. Unlike self-balancing binary search trees, the B-tree is optimized for systems that read and write large blocks of data.

R-tree

R-trees [14] are tree data structures used for spatial access methods, such as indexing multi-dimensional information like geographical coordinates. The key idea of the data structure is to group nearby objects and represent them with their minimum bounding rectangle in the next higher level of the tree; the "R" in R-tree is for rectangle. Since all objects lie within this bounding rectangle, a query that does not intersect the bounding rectangle also cannot intersect any of the contained objects. At the leaf level, each rectangle describes a single object; at higher levels the aggregation of an increasing number of objects. This can also be seen as an increasingly coarse approximation of the data set. R-tree's searching algorithms use the bounding boxes to decide whether or not to search inside a

subtree. In this way, most of the nodes in the tree are never read during a search. Like B-trees, this makes R-trees suitable for large data sets and databases, where nodes can be paged to memory when needed, and the whole tree cannot be kept in main memory.

GiST

GiST [16] stands for Generalized Search Tree. It is a balanced, tree-structured access method, that acts as a base template in which to implement arbitrary indexing schemes. B-trees, R-trees and many other indexing schemes can be implemented in GiST. GiST tree nodes contain a pair of values in the form of (p, ptr). The first value, p, is used as a search key and the second value, ptr, is a pointer to the data if the node is a leaf or a pointer to another node if the node is intermediate. The search key p represents an attribute that becomes true for all data that can be reached through the pointer ptr. Also, GiST index can handle any query predicate, as long as certain functions, that influence the behavior of the search keys, are implemented.

PostGIS can use GiST index at a table attribute. When a table column is a geography data type, then GiST will use an improved version of an R-tree index (R-tree-over-GiST scheme). PostgreSQL doesn't allow, at the most recent versions, the use of standard R-tree, because this type of index cannot handle attributes with size bigger than 8K and fail when a geography column is null.

2.3 HDFS

The Hadoop Distributed File System (HDFS) [22] is a distributed file system designed to run on commodity hardware. HDFS is highly fault-tolerant and is designed to be deployed on low-cost hardware. HDFS provides high throughput access to application data and is suitable for applications that have large data sets. HDFS is designed to reliably store very large files across machines in a large cluster. It stores each file as a sequence of blocks; all blocks in a file except the last block are the same size. The blocks of a file are replicated for fault tolerance.

The NameNode is the centerpiece of an HDFS file system. It keeps the directory tree of all files in the file system, and tracks where across the cluster the file data is kept. It does not store the data of these files itself. Client applications talk to the NameNode whenever they wish to locate a file, or when they want to add/copy/move/delete a file. The NameNode responds the successful requests by returning a list of relevant DataNode servers where the data lives. A Datanode stores data in the HDFS. Client applications can talk directly to a DataNode, once the NameNode has provided the location of the data. The DataNodes also perform block creation, deletion, and replication upon instruction from the NameNode.

2.4 HBase

HBase is an open source, non-relational (NoSQL), distributed database modeled after Google's BigTable [12]. It is developed as part of Apache Software Foundation's Apache Hadoop project and runs on top of HDFS.

2.4.1 Data model

HBase works best for sparse data sets, which are common in many Big Data use cases. Unlike relational database systems, HBase does not support a structured query language like SQL; in fact, HBase isn't a relational data store at all. HBase is a column-oriented database and follows a key/value data storage model. The main characteristics of the data model are the following:

- Table: HBase organizes data into tables as most database systems.

- **Row:** Within a table, data are stored according to its row. Rows are uniquely identified by a key. Row keys are treated as byte arrays, thus they do not have a specific data type.
- **Column Family:** Column families group data within a row, impacting their physical arrangement. They are stored together on disk, which is why HBase is referred to as a column-oriented data store. Column families must be defined up front, during table creation.
- **Column Qualifier:** Data within a column family are addressed via its column qualifier, or simply, column. Column qualifiers need not be consistent between rows. Like row keys, column qualifiers do not have a data type and are always treated as a byte array.
- **Cell:** A combination of row key, column family, and column qualifier uniquely identifies a cell. The data stored in a cell are referred to as that cell's value. Values also do not have a data type and are always treated as a byte array.
- **Timestamp:** Values within a cell are versioned. Versions are identified by their version number, which by default is the timestamp of when the cell was written. If a timestamp is not specified during a write, the current timestamp is used. If the timestamp is not specified for a read, the latest one is returned. The number of cell value versions retained by HBase is configured for each column family. The default number of cell versions is three.

For example, a table in HBase can be the following:

Each cell has multiple versions, typically represented by the timestamp of when they were inserted into the table

Timestamp1 Timestamp2

Row Key	Column Family - Personal			Column Family - Office	
	Name	Residence	Phone	Phone	Address
00001	John	415-111-1234	415-212-5544	1021 Market St	
00002	Paul	408-432-9922	415-212-5544	1021 Market St	
00003	Ron	415-993-2124	415-212-5544	1021 Market St	
00004	Rob	818-243-9988	408-998-4322	4455 Bird Ave	
00005	Carly	206-221-9128	408-998-4325	4455 Bird Ave	
00006	Scott	818-231-2566	650-443-2211	543 Dale Ave	

ie table is lexicographically sorted on the row keys

Cells

Figure 2.1: A table in HBase consisting of two column families

2.4.2 Architecture components

HBase has the following three major architecture components:

- **Region server:** Regions are parts of HBase tables, splitted accross the region servers of the distributed system and managed by a region server. Region server communicates with the client and handles data-related operations. Also, it handles read and write requests for all regions under it and decides the size of the regions.
- **Master server:** Master server assigns regions to the region servers and handles the load balancing of the regions across region servers. Moreover, it is responsible for schema changes and other metadata operations such as creation of tables and column families.
- **Zookeeper:** Zookeeper server is responsible to handle and synchronize the communication between the master and region servers, providing fault tolerance. More specifically, ZooKeeper is a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services.

2.4.3 Coprocessors

HBase coprocessors [6] enable distributed computation directly within the HBase server processes on the server's local data. The idea of HBase Coprocessors was inspired by Google's BigTable coprocessors. Coprocessors can be loaded globally on all tables and regions hosted by the region server, these are known as system coprocessors; or the administrator can specify which coprocessors should be loaded on all regions for a table on a per-table basis, these are known as table coprocessors. There are two categories of HBase coprocessors, the observer and the endpoint coprocessors. The idea behind observers is that we can insert user code by overriding upcall methods provided by the coprocessor framework. The callback functions are executed from core HBase code when certain events occur. Endpoints, on the other hand, are more powerful, resembling stored procedures. One can invoke an endpoint at any time from the client. The endpoint implementation will then be executed remotely at the target region or regions, and results from those executions will be returned to the client. In essence, coprocessors implement the idea of "move computation near the data", a method that Google heavily utilizes to increase performance: it is much cheaper to transfer a small piece of code to be executed in parallel in various servers instead of moving data where the code is installed.

The Endpoint is an interface for dynamic RPC extension. The endpoint implementation is installed on the server side and can then be invoked with an HBase RPC. The client library provides convenience methods for invoking such dynamic interfaces. The application code client side performs a batch call. This initiates parallel RPC invocations of the registered dynamic protocol on every target table region. The results of those invocations are returned as they become available. The client library manages this parallel communication on behalf of the application, messy details such as dealing with retries and errors, until all results are returned (or in the event of an unrecoverable error). Then the client library rolls up the responses into a Map and hands it over to the application. If an unrecoverable error occurs, then an exception will be thrown for the application code to catch and take action.

2.5 Google Maps Directions API

The Google Directions API [2] is a service that calculates directions between locations using an HTTP request. A Directions API request takes the following form:

`http://maps.googleapis.com/maps/api/directions/output?parameters`

where output determines the type of file that will contain the answer to the HTTP request. The type of the output file can be either json or xml. The field parameters defines the parameters and their values that will come along with the HTTP request. Certain parameters are required while others are optional. As is standard in URLs, all parameters are separated using the ampersand (&) character. Required parameters are the origin and destination locations between which we request directions through the Directions API. These locations can be in the form of textual geographical coordinates as a value pair of latitude and longitude. Location can also be inserted as the address of the specific point. In this case, Directions service will geocode the string and convert it to a latitude/longitude coordinate in order to calculate directions. Other optional parameters, that can be defined, are the mode of transportation, which can be driving, walking, bicycling, or transit. For example, a complete HTTP request to the Directions API can be the following:

`http://maps.googleapis.com/maps/api/directions/json?origin=37.976159,
23.776274&destination=37.978180, 23.768957&mode=walking`

2.5.1 JSON response file

JSON (JavaScript Object Notation) file format, is an open standard format that uses human-readable text to transmit data objects consisting of attribute-value pairs. Directions responses in json type format contain the following fields:

Status

The status field within the Directions response object contains the status of the request, and may contain debugging information to help you track down why the Directions service failed. The status field may contain the following values:

- OK, the response contains a valid result.
- NOT_FOUND, at least one of the locations specified in the request's origin, destination, or waypoints could not be geocoded.
- ZERO_RESULTS, no route could be found between the origin and destination.
- INVALID_REQUEST, the provided request was invalid. Common causes of this status include an invalid parameter or parameter value.
- REQUEST_DENIED, the service denied use of the directions service by your application.
- UNKNOWN_ERROR, a directions request could not be processed due to a server error.
- OVER_QUERY_LIMIT, the service has received too many requests from your application within the allowed time period. More specifically, users of the free Directions API are able to perform at most 2 HTTP request per second and 2500 HTTP requests per day to the API.

For example:

```
"status" : "OK"
```

Routes

Routes field is a JSON array, whose elements represent a different route between the origin and destination locations. Each route can contain the following fields:

- summary: a short textual description for the route, suitable for naming and disambiguating the route from alternatives.
- legs[]: an array which contains information about a leg (i.e. a part) of the route, between two locations within the given route. If there are no waypoints defined, the array will include one element, which will be the total route.
- waypoint_order: an array indicating the order of any waypoints in the calculated route.
- overview_polyline: a single points object that holds an encoded polyline representation of the route. This polyline is an approximate (smoothed) path of the resulting directions.
- bounds: the viewport bounding box of the overview_polyline.
- copyrights: the copyrights text to be displayed for this route.
- warnings[]: an array of warnings to be displayed when showing these directions.

For example:

```
"bounds" : {
  "northeast" : {
    "lat" : 37.978186099999999,
    "lng" : 23.7762659
  },
  "southwest" : {
    "lat" : 37.9761428,
```

```

    "lng" : 23.7689148
  }
},
"copyrights" : "Map data ©2015 Google",
"legs" : [
  ... ]
"overview_polyline" : {
  "points" : "{exfFuxbpCg@fCkbpK}B-MkB|KIf@QG"
},
"summary" : "Λεωφ. Στρατάρχου Αλεξάνδρου Παπάγου ",
"warnings" : [
  "Walking directions are in beta. Use caution - This route may be missing
    sidewalks or pedestrian paths."
],
"waypoint_order" : []

```

legs[]

Each element in the legs array specifies a single leg of the journey from the origin to the destination in the calculated route. For routes that contain no waypoints, the route will consist of a single "leg," but for routes that define one or more waypoints, the route will consist of one or more legs, corresponding to the specific legs of the journey.

Each leg within the legs field(s) may contain the following fields:

- **steps[]**: an array of steps denoting information about each separate step of the leg of the journey.
- **distance**: the total distance covered by this leg, where value field contains the distance in meters and text field contains a human-readable representation of the distance, displayed in units as used at the origin
- **duration**: the total duration of this leg, where value field contains the duration in seconds and text field contains a human-readable representation of the duration.
- **start_location**: contains the latitude/longitude coordinates of the origin of this leg. Because the Directions API calculates directions between locations by using the nearest transportation option (usually a road) at the start and end points, start_location may be different than the provided origin of this leg if, for example, a road is not near the origin.
- **start_address**: contains the human-readable address (typically a street address) reflecting the start_location of this leg.
- **end_location**: contains the latitude/longitude coordinates of the given destination of this leg. Because the Directions API calculates directions between locations by using the nearest transportation option (usually a road) at the start and end points, end_location may be different than the provided destination of this leg if, for example, a road is not near the destination.
- **end_address**: contains the human-readable address (typically a street address) reflecting the end_location of this leg.

For example:

```

  "distance" : {
    "text" : "0.7 km",
    "value" : 690
  },
  "duration" : {
    "text" : "7 mins",
    "value" : 434
  },

```

steps[]

- `html_instructions`: formatted instructions for this step, presented as an HTML text string.
- `distance`: the distance covered by this step until the next step.
- `duration`: the time required to perform the step, until the next step.
- `start_location`: the geographical coordinated of the location of the starting point of this step.
- `end_location`: the geographical coordinated of the location of the last point of this step.
- `polyline`: a single points object that holds an encoded polyline representation of the step. This polyline is an approximate (smoothed) path of the step.
- `travel_mode`: the mode of transportation, as it is defined in the mode parameter of the HTTP request.

```
"distance" : {  
    "text" : "0.7 km",  
    "value" : 680  
},  
"duration" : {  
    "text" : "7 mins",  
    "value" : 427  
},  
"end_location" : {  
    "lat" : 37.9781038,  
    "lng" : 23.7689148  
},  
"html_instructions" : "Head \u003cb\u003ewest\u003c/b\u003e on \u003cb\u003eΛεωφου\u003e.  
Στρατάρχου Αλεξάνδρου Παπάγου \u003c/b\u003e toward \u003cb\u003eΓρυ\u003e. Κουσιδη\u003c/b\u003e",  
"polyline" : {  
    "points" : "{\xfFuxbpCQv@ETCLKj@CPYxAi@|Cc@fCm@lDg@pCCNc@nC_AnFCHYhBMx@if@"  
},  
"start_location" : {  
    "lat" : 37.9761428,  
    "lng" : 23.7762659  
},  
"travel mode" : "WALKING"
```

2.5.2 Encoded Polyline Algorithm Format

The JSON response file contains an encoded polyline as a representation of the route. This representation is the result of encoding the geographical coordinates of the intermediate points of the route. Google Maps API follows a specific encoding algorithm [3] of the latitude and longitude of any point. The steps for encoding such a signed value indicating the latitude or longitude are specified below:

1. Multiply the initial signed value by $1e5$, rounding the result.
2. Convert the result to its binary equivalent value, using its two complement.
3. Left-shift the binary value one bit.
4. If the original decimal value is negative, invert this encoding.
5. Break the binary value out into 5-bit chunks (starting from the right hand side).
6. Place the 5-bit chunks into reverse order.
7. OR each value with $0x20$ if another bit chunk follows.
8. Convert each value to decimal.
9. Add 63 to each value.
10. Convert each value to its ASCII equivalent.

Following these steps in reverse we can decode an encoded polyline into the appropriate list of geographical coordinates and extract the GPS traces that define the route.

2.6 Google Maps Static Maps API

The Google Static Map service [4] creates a map based on URL parameters sent through a standard HTTP request and returns the map as an image that can be displayed on a web page, or accessed through a URL. A simple HTTP request to the Google Static Maps API takes the following form:

`https://maps.googleapis.com/maps/api/staticmap?parameters`

The Static Maps API defines map images using the following URL parameters:

Location Parameters

- center: defines the center of the map, equidistant from all edges of the map. This parameter takes a location as either a comma-separated latitude,longitude pair or a string address identifying a unique location on the face of the earth.
- zoom: defines the zoom level of the map, which determines the magnification level of the map.

Map Parameters

- size: the rectangular dimensions of the map image. This parameter takes a string of the form {horizontal_value}x{vertical_value}.
- scale: affects the number of pixels that are returned. scale=2 returns twice as many pixels as scale=1 while retaining the same coverage area and level of detail.
- format: the format of the resulting image. Possible formats include PNG, GIF and JPEG types.

- **maptype:** the type of map to construct. There are several possible maptype values, including roadmap, satellite, hybrid, and terrain.
- **language:** the language to use for display of labels on map tiles.

Feature Parameters

- **markers:** one or more markers to attach to the image at specified locations. This parameter takes a single marker definition with parameters separated by the pipe character (|).
- **path:** a single path of two or more connected points to overlay on the image at specified locations. This parameter takes a string of point definitions separated by the pipe character (|). In addition, the encoded polyline representation can be used instead of multiple points, as long as the prefix **enc:** is used in the value of the parameter path.
- **visible:** one or more locations that should remain visible on the map, though no markers or other indicators will be displayed.

For example, a HTTP request at Google Static Maps API can be the following:

```
https://maps.googleapis.com/maps/api/staticmap?&size=1000x1000&markers=label:A|17.7466,-
64.703&markers=label:C|17.7453,-
64.7019&path=color:blue|enc:gcikBvh|iKFJ??nAaAFGjA}@rAiAJK
```

2.7 Mathematical Background

2.7.1 Random Variables

A random variable X is a measurable function from the set of possible outcomes Ω , $X: \Omega \rightarrow A$ or $X: \Omega \rightarrow \mathbb{R}$, where A is a subset of the set of real numbers. If A is a set of discrete values, such as the set of integer numbers, then the random variable is discrete. If A is a set of infinite possible values, then the random variable is continuous. In general, a random variable is a variable whose value is subject to variations due to chance. A random variable can take on a set of possible different values (similarly to other mathematical variables), each with an associated probability if the variable is discrete, or with a probability density function if the variable is continuous. The function of the random variable is called probability distribution.

2.7.2 Uniform Distribution

The discrete uniform distribution is a symmetric probability distribution whereby a finite number of values are equally likely to be observed; every one of n values has equal probability $1/n$.

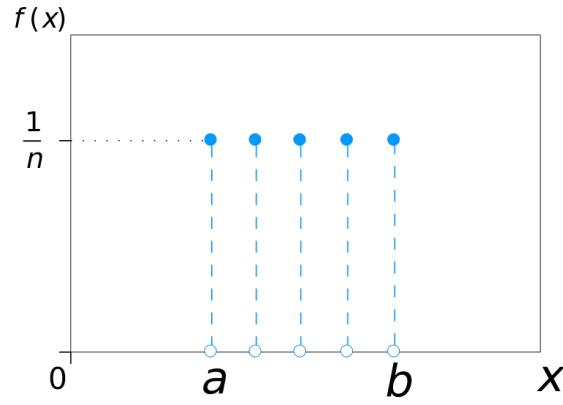


Figure 2.2: Discrete uniform probability mass function

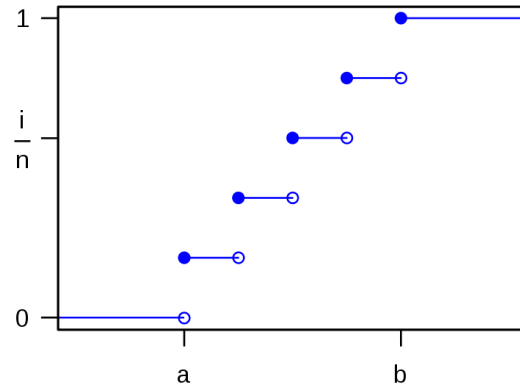


Figure 2.3: Discrete uniform cumulative distribution function

2.7.3 Normal Distribution

Normal distribution, known also as Gauss distribution, refers to continuous random variables. The normal distribution is remarkably useful because of the central limit theorem. In its most general form, it states that averages of random variables drawn from independent distributions are normally distributed. The probability density of the normal distribution is:

$$P(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

where μ is the mean of the distribution and σ is the standard deviation of the distribution.

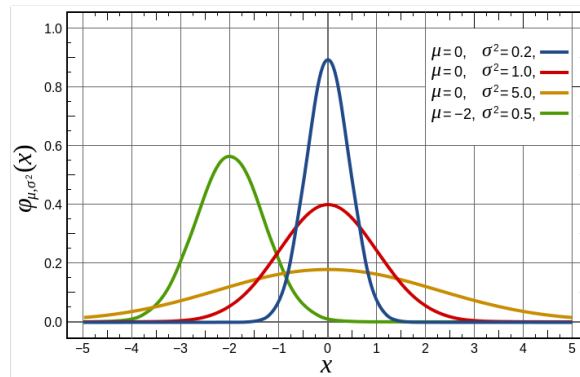


Figure 2.4: Probability density function

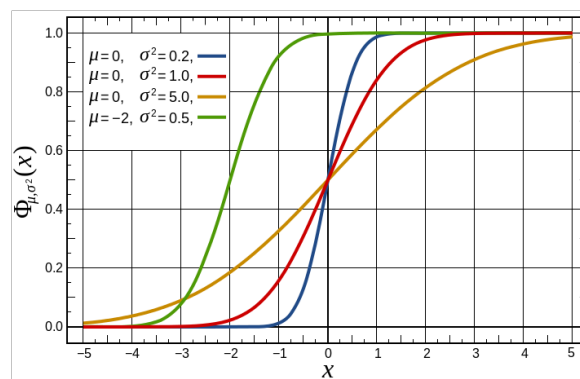


Figure 2.5: Cumulative distribution function

According to an empirical rule, about 68% of values drawn from a normal distribution are within one standard deviation σ away from the mean; about 95% of the values lie within two standard deviations; and about 99.7% are within three standard deviations. This fact is known as the 68-95-99.7 rule, or the 3-sigma rule.

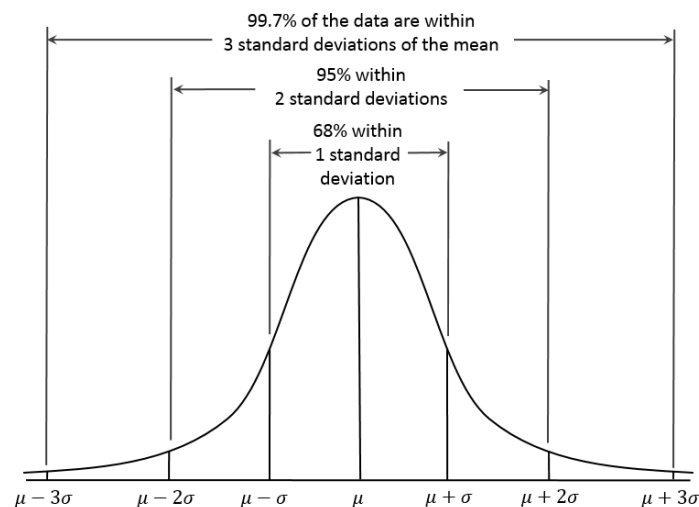


Figure 2.6: 3-sigma empirical rule for the normal distribution

2.7.4 Bernoulli Distribution

The Bernoulli distribution is the probability distribution of a random variable which takes value 1 with success probability p and value 0 with failure probability $q=1-p$. It can be used, for example, to represent the toss of a (not necessarily fair) coin, where "1" is defined to mean "heads" and "0" is defined to mean "tails" (or vice versa).

If X is a random variable with this distribution, we have:

$$P(X = 1) = p$$

$$P(X = 0) = 1 - p$$

The experiment is called fair if $p=0.5$ and can represent the toss of a fair coin.

Chapter 3

Generator

3.1 Storage of source data

The generator of spatio-temporal and textual data uses real point of interests on the map as a source of data. These data were extracted from the online travel service TripAdvisor, using an already implemented crawler [17]. These points of interest are identified on the map by their geographical coordinates, as well as their address. Also, they come with ratings and reviews by real TripAdvisor users. The number of such available points is 136409, as they were extracted by a 13GB response json file from the crawler. The points were stored in PostgreSQL. The database schema used, contains two tables, one for the points of interest and one for the reviews and ratings of these points and is the following:

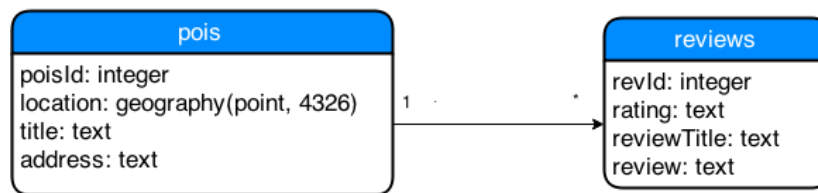


Figure 3.1: Database schema for source data

More specifically, the table for the points of interest has the following attributes:

- `poisId`: identifying number for the point of interest, primary key.
- `location`: geographical coordinates of the point. Usage of the geographic data type introduced by the extension PostGIS of PostgreSQL.
- `title`: the name of the point of interest.
- `address`: the address of the point on the map.

A point of interest can have many reviews, issued by different users, thus the two tables are associated by 1-to-many. The reviews table contains the following attributes:

- `revId`: the identifying number of the point of interest that the review refers to, foreign key.
- `rating`: rating of the point in a scale of 1 to 5.
- `reviewTitle`: the title of the review.
- `review`: the text of the review.

After storing all points of interest into the database, we create an B-tree index to the attribute `poisId` of the table **pois** and to the attribute `revId` of the table **reviews**, as well. Additionally, we create a GiST index to the attribute `location` of the table **pois**. As `location` has a geographic data type, GiST will implement an improved R-tree index, as described in section 2.2.2. In this way, searching points of interest according to their location or identifying number, can be done fast and efficiently.

3.2 Generator attributes

The following class diagram indicates the main attributes used in the description of the generator's design.

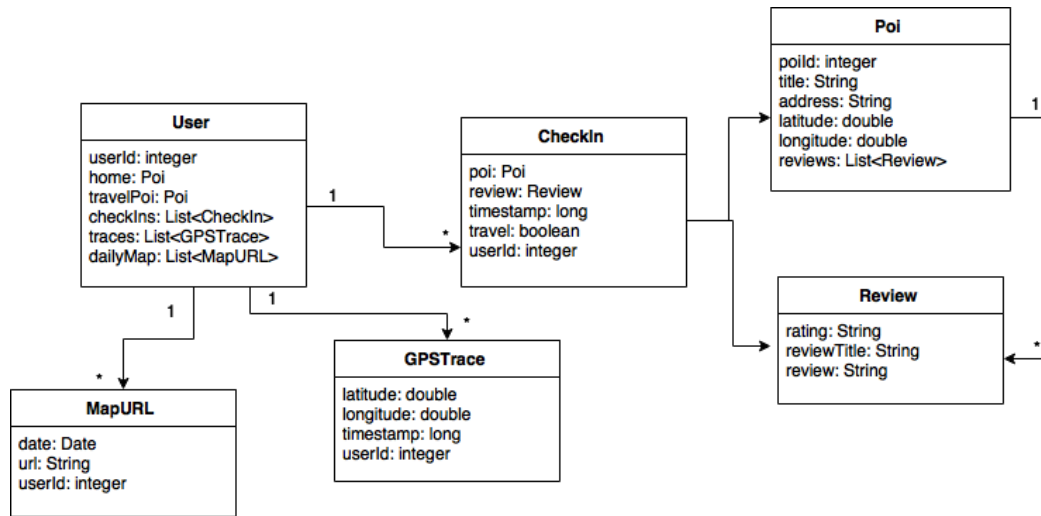


Figure 3.2: Class diagram of generator's attributes

3.2.1 User

The user produced by the generator has the following fields:

- **userId**: user's identifier.
- **home**: user's home defined as Poi.
- **travelPoi**: the central location of the current trip, defined as Poi.
- **checkIns**: list of user's check-ins during the days for which the generator produced data.
- **traces**: list of GPS traces indicating the user's daily routes during the days for which the generator produced data.
- **dailyMap**: list of URLs to static maps images showing the daily user routes.

3.2.2 Check-in

A check-in to a point of interest (POI) has the following fields:

- **poi**: the point of interest (POI) where the check-in was made.
- **review**: the review that the user made for the specific POI on the check-in.
- **timestamp**: the time when the check-in was made, in representation of type long.
- **travel**: boolean value indicating whether the check-in was made during a user's trip.
- **userId**: the user's identifier.

3.2.3 GPS trace

A user's route consists of multiple GPS traces, which illustrate the path that the user followed. A GPS trace contains the following information:

- (latitude, longitude): the geographical coordinates of the point on the map where the GPS trace was taken.
- timestamp: the time when the GPS trace was taken.
- userId: user's identifier whose route contains the specific GPS trace.

3.2.4 Static route map - MapURL

A user's daily route is depicted on a static map, using Google Static Maps API. The map shows the points of interest that the user visited using markers, and the exact path from one POI to another using blue continuous lines. The image of the map can be accessed through a respective URL.

- url: the URL that directs to the image of the map.
- date: the date when the user walked the route showed on the map.
- userId: the identifier of the user who walked the specific route that date.

3.2.5 Point of interest - Poi

A point of interest has the following attributes:

- poiId: point's identifier.
- title: point's name.
- address: point's address.
- (latitude, longitude): point's geographical coordinates.
- reviews: a list with the available reviews for the specific point.

3.2.6 Review

A review to a point of interest has the following fields:

- rating: rating of the POI in scale 1 to 5.
- reviewTitle: the title of the review.
- review: the text of the review.

3.3 Input parameters

The generator takes as input the next parameters:

- userIdStart: identifier of the first user for whom the generator will create daily routes.
- userIdEnd: Respectively, the identifier of the last user created.
- chkNumMean: The mean of the number of daily check-ins to points of interest, which will follow a normal distribution.

- **chkNumStDev**: The standard deviation of the normal distribution that the number of daily check-ins follow.
- **chkDurMean**: The mean of the duration that each visit will last, which will follow a normal distribution.
- **chkDurStDev**: The standard deviation of the normal distribution of the duration of user's visits.
- **dist**: The maximum distance in meters in which a user can walk from one point of interest to the next one.
- **maxDist**: The maximum distance in meters from a user's home, in which a user can walk every day.
- **startTime**: The time when the first check-in of the day will occur.
- **endTime**: The time when the last check-in of the day will occur.
- **startDate**: The first day starting from which the generator will create daily routes.
- **endDate**: Respectively, the last day that the generator will create daily routes.
- **outCheckIns**: The output file storing the daily check-ins of all users created.
- **outTraces**: The output files storing the GPS traces of the daily routes of all users created.
- **outMaps**: The output file storing the URLs for the daily maps depicting the daily routes of all users created.

3.4 Generator's configuration parameters

The generator's main function is to create user's daily routes, which will contain information simulating real social media data. These daily routes consist of visits (check-ins) to points of interest and paths that the user followed from one point to the next one. The generator takes as input the values of the corresponding parameters and based on the stored source data, creates daily routes for $(\text{userIdEnd} - \text{userIdStart} + 1)$ users between the start and end dates. The decisions that the generator makes, such as which and how many points the user will visit a specific day, or the duration of each visit, are made using random factors. In this way, the produced data can be as realistic as possible. The generator stores the created data in three files, one for all the user check-ins (**outCheckIns**), one for the all the GPS traces (**outTraces**) and one for the all the daily maps (**outMaps**).

Briefly, the daily routes that the generator produces are designed in such a way that allows each user to visit places by walking from one point to the other. Each day, the user leaves his home and visits certain places that are in walking distance. The number of visits per day is decided in a random way. Moreover, the user stays for a couple of hours, also defined in a random way, at a point of interest and leaves a rating and review for that specific POI. Moving on, he walks to the next point, following the path that Google Directions API provides, and checks in to the next place. The generator enables a user to travel, following the same daily route strategy in the location of his trip.

3.4.1 Home

The location of a user's home is the center around which he walks every day. The points of interest that a user visits are selected in a random way, thus there can exist routes that don't make sense. For example, a user can walk one day in Greece, the next day in France and the other day again back in Greece. This is the reason why the location of a user's home is specified, so that a user can walk in a sensible distance from his home every day. More specifically, a user can walk in a

range of maxDist meters from his home, as those are defined by the according input parameter. In the implementation of the generator, the home of a user is defined as the point where his first ever visit is made on startDate. The choice of this point is done using a generator of random numbers which follow a uniform distribution. The range of the distribution are the total number of source POIs stored in PostgreSQL.

3.4.2 Trip

A user created by the generator is capable to travel, so as to be able to visit places that are further than maxDist meters from his home. In this way, routes one day in Greece and the next one in France make sense. However, a central point around which the user will walk during his trip has to be defined. In the implementation of the generator, this point is chosen to be the first ever point that the user visits during the current trip. The choice of that point is done using a generator of random number that follow a uniform distribution. The range of the distribution is the number of available source POIs.

As far as the duration of each trip is concerned, each user can travel for a time interval that equals the 10% of the total time interval for which the generator produces daily routes. The days of the time interval can be spread out to multiple trips. The duration of each trip is defined using a generator of random numbers which follow a normal distribution. The mean of this distribution is declared to be 5 and the standard deviation is 2. In this way, according to the 3-sigma empirical rule for the normal distribution, 95% of the random trip durations will be between 1 and 9 days, which is reasonable for short and long trips. To sum up, if the user is about to travel the next days, the duration of the trip is defined in a random way and if the duration of the trip doesn't exceed the available travel days, the prospective trip begins. If the duration of the trip exceeds the available days, then the trip is calculated to last for the available days.

Finally, the decision whether the user will begin a trip or not is made in a random way, using a generator of random numbers that follow the Bernoulli distribution. Thus, the generator decides every day whether or not to start a trip for the current user. This decision resembles a fair coin toss. Therefore, if the generator decides to start a trip for the current user, then it decides, in a random way as well, the duration and the location of the trip. Each daily route during the trip has to be in maxDist range from the trip's central location.

3.4.3 Check-in POI

The number of daily check-ins is defined using a generator of random numbers that follow a normal distribution with a mean value determined by the input parameter chkNumMean and standard deviation determined by the input parameter chkNumStDev. Therefore, every day the generator picks a different number of daily check-ins, which according to the 3-sigma empirical rule for the normal distribution, will most probably be around the mean value.

The locations of check-ins are determined in a random way by the generator. More specifically, the choice of the point of interest, where the first check-in of the day will take place, is made using a generator of random number who follow a uniform distribution. The range of this distribution is the number of available POIs who are located in maxDist range from the user's home. In this way, the generator will pick a random POI between those who are in walking distance from the user's home. If the user travels, then the generator will choose a random POI between those who are in maxDist range for the trip's central location. In order to find those points and select them from the PostgreSQL database, the function ST_DWithin function is used, which is available through the PostGIS extension. This function returns true for the points which are in the desired distance from the user's home, calculating the distance using the geographical coordinates of the points. The search of these points in the PostgreSQL table is donw efficiently, due to the GiST index. The generator will assemble all the points that are in the desired range, and choose a random one as the first point that the user visits in the current day.

Using the same strategy, the generator chooses all next points of interest that the user will visit the specific day. However, the next points visited should be in a smaller distance from the first POI visited. This distance is defined by the input parameter `dist`. Also, a user is not allowed to visit the same place twice during the day, so every next check-in should be in a POI not visited that specific day.

3.4.4 Check-in review

The generator assigns a rating and review to every user's check-in. More specifically, every source point of interest as it is stored in PostgreSQL database, contains certain reviews for the specific point. The generator chooses randomly a review amongst the available for the POI, using a generator of random numbers that follow a uniform distribution. The range of the distribution is the number of available reviews for the specific POI.

3.4.5 Path between check-ins

The generator issues an http request to Google Directions API, in order to obtain information about the path that a user will follow in order to walk from one point of interest to the next one. Since the source data contain the geographical coordinates of every point of interest, the generator has access to the latitude and longitude of every available point. Thus, he sets as value of the origin parameter to the http request the geographical coordinates of the current point where the user checked in and as destination the coordinates of the next point of interest that the user will visit, as it was selected in a random way. Also, he specifies the parameter mode into walking, because the user walks from one POI to the next one. Finally, the response file of the request will be in json file format. For example, a request to Google Directions API can be the following:

```
http://maps.googleapis.com/maps/api/directions/json?origin=37.976159,  
23.776274&destination=37.978180, 23.768957&mode=walking
```

The json response file contains information about the path on the map, that the user will have to follow in order to get to his destination, as well as the duration of his walk to the destination. We extract from the json file, the field `polyline` from each step of the route. The `polyline` holds an encoded representation of the step's path. We decode each step's `polyline` using the reverse Encoded Polyline Algorithm Format, as it was described in section 2.4.2. In this way, we have access to a list of geographical coordinates indicating all the points on the map, that the user will walk through from the origin to the destination. These points will be stored as GPS traces, representing the user's route. For each path, the starting and ending points, which are the points of interest, will also be stored as GPS traces.

3.4.6 Timestamp

As far as time is concerned, the generator defines that the first check-in of the day will happen at the time defined by the input parameter `startTime`. The duration of the visit to a point of interest is set using a generator of random numbers that follow a normal distribution. The mean value of the distribution is defined by the `chkDurMean` input parameter, and the standard deviation by the `chkDurStDev` parameter. The time when the next check-in will occur is set to be the time when the previous one happened, plus the duration of the previous visit and the duration of the walk from the previous point to the next one. The duration of the walk, is extracted from the field `duration` of the Google Directions json response file. If the time that the next check-in will happen exceeds the time defined by the input parameter `endTime`, then the next check-in won't occur, and the check-ins end at that time for that specific day. Each check-in is timestamped using the Unix Time Stamp, which is a long integer representation of the date and specific time (UTC timezone) of the event.

Concerning the timestamp of the GPS traces, they are calculated through the json response file. The duration of the route from the origin to the destination is splitted up by the number of points decoded from the polyline. Therefore, the timestamp of the first GPS trace of the path is set to be the time the visit at the origin ended plus the fraction of the divided time needed to get to that point on the map.

3.4.7 Static map

The generator uses a static map in order to depict the user's daily route, by issuing an http request to the Google Static Maps API. The points where a user checked in during the day are distinctly visible on the map with markers. These markers are also named using capital letters of the alphabet in order to show the order that the user visited them. The generator uses the stored polylines, as they were extracted by the json response file, in order to define the path that the map will indicate using a blue continuous line. The image of the map is accessible through the URL that forms the http request to the API. The generator stores the URLs created at the corresponding output file (outMaps), so that the image of the map can be accessible by the users of the generator.

For example, a request to Google Static Maps API, as it is created by the generator, can be the following:

```
https://maps.googleapis.com/maps/api/staticmap?&size=1000x1000&markers=label:
A|44.7698, -69.7215&markers=label:B|44.7651, -69.7189&markers=label:
C|44.7639, -69.7196&markers=label:D|44.7656, -69.717&path=color:blue|enc:
gbgpGjnphL@FZKtE_BxEeB'Bk@z@[bA]tBo@HG~Ai@JMFG?A@A?A@C?C?E?C?EAE?I??M_??L^?
?Tj@Vl@LNDFFDDFDDDB@DBD@HBLBRBB?H@JAF?DAFAHCFDCBCDEGDU??OTEFEDCBEBGBIBG@E@G?K@
IAC?SCMCICEAECCAEEGEEEGMO??Wm@Uk@cAyCs@_C??bA_A
```

The map that corresponds to the above URL is:

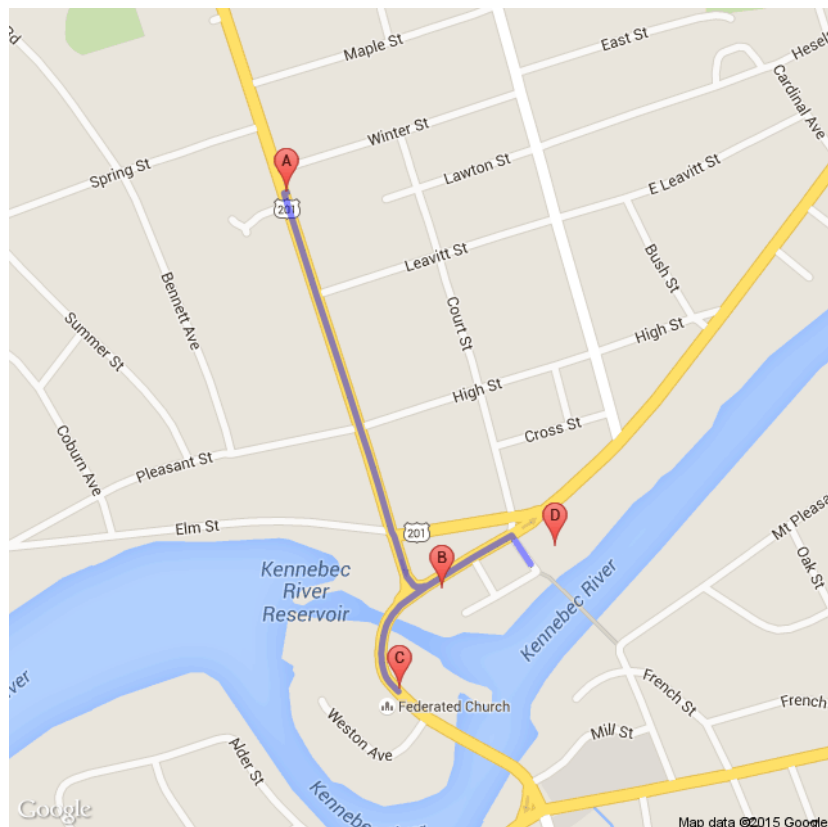


Figure 3.3: Example of static map image

Chapter 4

Dataset

4.1 Generator input parameters

We were able to create a large dataset using the generator. More specifically, we set the input parameters of the generator to the following values:

$$\text{chkNumMean} = 5 \text{ } \text{chkNumStDev} = 2$$

The number of daily check-ins will be random and 95% of these random values will be between 1 and 9, according to the 3-sigma empirical rule for the normal distribution.

$$\text{chkDurMean} = 2 \text{ } \text{chkDurStDev} = 0.1$$

The duration of each visit will be random and 95% of these random values will be between 1.8 and 2.2 hours.

$$\text{maxDist} = 50000.0 \text{ } \text{dist} = 500.0$$

Each user will be able to visit places that are in 50000 meters radius from his home or travel location. Also, he is allowed to walk in a 500 meters radius from the one place he visits to the next one.

$$\text{startTime} = 9 \text{ } \text{endTime} = 23$$

Each user will visit the first place of the day at 9 am. Moreover, the last daily check-in should take place no later than 11 pm.

$$\text{startDate} = 01-01-2015 \text{ } \text{endDate} = 03-01-2015$$

The generator will produce daily routes for the time period between the 1st of January 2015 and the 1st of March 2015.

Therefore, the generator run with the specific input parameters in order to create a large dataset of users and their respective daily routes.

4.2 Generator deployment setup

Respecting the request restrictions for the users of the free Google Directions API to 2500 requests per day to the API, the generator could run once a day creating 14 users at a time for the specific input parameters. Therefore, we set up a cluster of 31 Virtual Machines (VMs) in order to be able to create a much bigger number of users per day. Each VM of the cluster had 1 CPU, 1 GB RAM and 10 GB disk. The cluster used resources located at the Computing Systems Laboratory at NTUA. The PostgreSQL database, where the source data were stored, was set up in a different VM where the generator run as

well. This specific VM had 2 CPU, 4 GB RAM and 40 GB disk in order to be able to store the source data. When the generator was running on the cluster VMs, a remote connection to the PostgreSQL database was established in order to gain access to the source data as well.

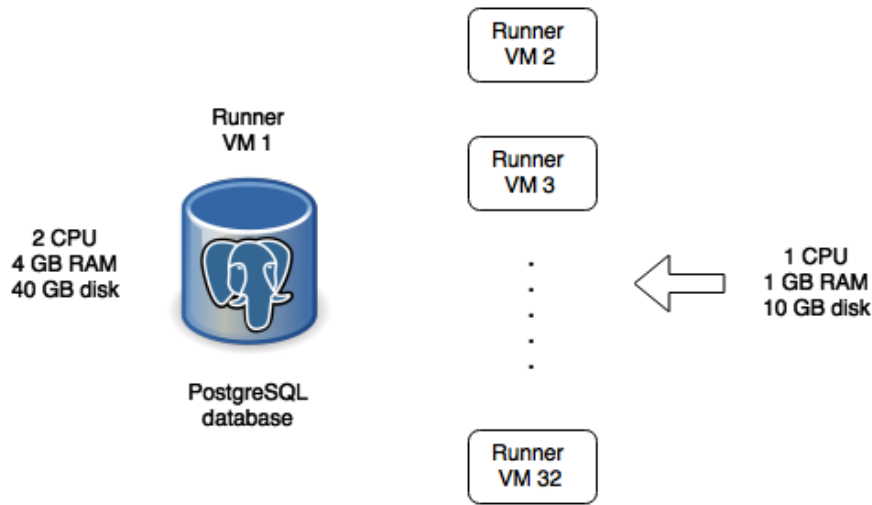


Figure 4.1: Generator execution architecture

In this way, we were able to run the generator on each VM collecting 448 users per day, creating 14 users per VM with the specified input parameters. At the end of the generator’s run period, we were able to have 9464 users and their 2 months daily routes. The created dataset was available in the generator’s two output files, one that stored the user’s daily check-ins and another one having the user’s total gps traces indicating the daily trajectories. More specifically, the generator created the following dataset:

Users	Check-ins	GPS traces
9464	1586537	38800019
3GB	641 MB	2.4 GB

Table 4.1: Dataset produced by the generator

4.3 Dataset storage model

The dataset created by the generator consists of 641 MB data of check-ins and 2.4 GB data of GPS traces. We also had available a 14 GB friend graph which we adjusted in order to match the number of 9464 users created by the generator. The overall dataset of check-ins, GPS traces and friends was stored into an HBase distributed database using the following data model.

4.3.1 Friends Table

The data about friends consist of user ids. For example, user no.1 has friends the users no.2, no.3 etc. This information is parsed from the friend graph and inserted into HBase in byte form using the following schema:

- Row: each row holds all the users that are friends of the user whose id is the key of the row.
- Column Family: there is one column family ‘friends’ including all the friends of one user.

- Column Qualifier: the qualifier that represents a single friend is the friend's user id.

Row	Friends		
Key	user id #1	user id #2	user id #3
1	145	2901	1204
2	3423	1023	965

Table 4.2: Example of HBase Table 'friends'

4.3.2 Check-ins Table

The user's check-ins are parsed from the corresponding generator's output file and inserted to HBase in byte form using the following schem:

- Row: each row holds all the check-ins of the user whose user id is the row key.
- Column Family: there is one column family 'checkIns' which holds all the check-ins of each user.
- Column Qualifier: the qualifier that represents a single user check-in is its timestamp.

Row	checkIns		
Key	timestamp #1	timestamp #2	timestamp #3
1	1420117200700	1420124471700	1420189200700

Table 4.3: Example of HBase Table 'check-ins'

4.3.3 GPS traces Table

The user's GPS traces are parsed from the corresponding generator's output file and inserted to HBase in byte form using the following schem:

- Row: each row holds all the GPS traces of the user whose user id is the row key.
- Column Family: there is one column family 'gpsTraces' which holds all the GPS traces of each user.
- Column Qualifier: the qualifier that represents a single user GPS trace is a string of the geographical coordinates combined with the timestamp of the GPS trace.

Row	gpsTraces	
Key	qualifier #1	qualifier #2
1	”(25.694, 79.836) 1420117200700”	”(25.69274, 79.8394) 1420124471700”

Table 4.4: Example of HBase Table 'gps-traces'

4.4 Dataset storage architecture

The overall dataset was stored into an HBase distributed database. We used the version 0.94.27 of HBase and the version 1.2.1 of Hadoop in order to utilize the HDFS. The HBase was set up over HDFS on a cluster of 32 VMs, consisting of 1 master and 32 region servers and 1 namenode and 32 datanodes. All different types of data were splitted into 32 parts, so that they are distributed equally into the region servers when the following tables were created. More specifically, when we created the table of 'friends' in HBase, we predefined the keys where the total table would be splitted into the region servers. Thus, the table was pre-splitted into 32 regions so that the data were equally divided into the region servers. The split into 32 parts was also done for the tables of 'check-ins' and 'gps-traces'. The master VM which contains the HBase master as well as the namenode has 2 CPU, 4GB RAM and 40 GB disk. The master is at the same time a region server and a datanode. The other 31 VMs holding the rest region servers and datanodes have 1 CPU, 2 GB RAM and 10 GB disk. The cluster is accomodated on a OpenStack managed cloud service located at the Computing Systems Laboratory in NTUA.

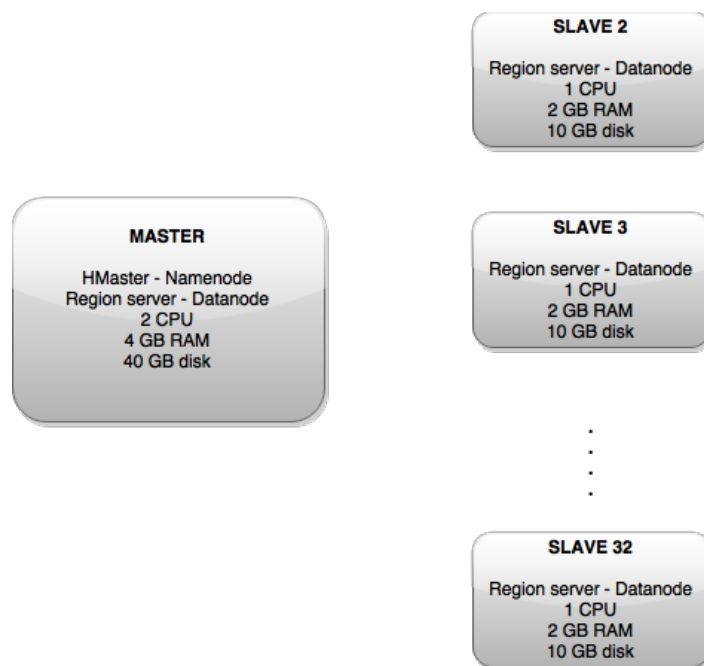


Figure 4.2: HBase cluster architecture

Chapter 5

Queries

After the insertion of all available data into the HBase cluster, we implemented several queries over the tables of 'friends' and 'check-ins'. These queries can be imposed to any social networking service that contains data about users that check in to several locations and have as friends other users of the service. As HBase is a NoSQL database and doesn't have a query execution language like SQL for example, the implementation of the queries was done using HBase coprocessors. In this way the computation of intermediate results and other complex calculations is transferred to the region servers that contain the respective data, decongesting the client from a heavy computational load.

5.1 Query implementation

We implemented the following queries over the data of tables 'friends' and 'check-ins' using HBase coprocessors.

5.1.1 Most visited POI query

This query contains the following question to the available data:

"Get the most visited points of interest of a certain user's friends"

The query is implemented in the following steps:

1. The client calls the coprocessor that returns the friends of the specified user. The coprocessor runs on the region server that contains the row of the table 'friends' that has as key the user id of the desired user. The client receives back a list of user id's that represent his friends.
2. The client splits the list of friends into sections according to the initial split of the 'check-ins' table into 32 regions. In this way, each splitted friends list will contain row keys that belong to only one region server.
3. The client issues a call to the coprocessor that calculates the most visited POIs for every splitted friend list. More specifically, the client starts a new thread that is responsible for calling the coprocessor and getting the result back. In this way, the client issues parallel calls to the coprocessor and the calculations to the respective region servers are done simultaneously.
4. The coprocessor that runs on each region server gets from the 'check-ins' table the rows that include all the check-ins of the user friends that are assigned to this region. He then iterates over each row in order to store and count how many times the user's friend has visited each POI. The POIs are stored into a hash table for faster calculations. At the end, he iterates over the hash table in order to get the POI that the specific friend visited the most times.
5. Finally, the client merges the results that he got back from each coprocessor call and forms the final result, which contains the place and how many times each of his friends went to their most visited POI.

5.1.2 News Feed query

This query contains the following question to the available data:

”Get the check-ins of all the friends of a specific user for a certain day into chronological order”

The query is implemented in the following steps:

1. The client calls the coprocessor that returns the friends of the specified user. The coprocessor runs on the region server that contains the row of the table 'friends' that has as key the user id of the desired user. The client receives back a list of user id's that represent his friends.
2. The client splits the list of friends into sections according to the initial split of the 'check-ins' table into 32 regions. In this way, each splitted friends list will contain row keys that belong to only one region server.
3. The client issues a call to the coprocessor that calculates the news feed for every splitted friend list. More specifically, the client starts a new thread that is responsible for calling the coprocessor and getting the result back. In this way, the client issues parallel calls to the coprocessor and the calculations to the respective region servers are done simultaneously.
4. The coprocessor that runs on each region server gets from the 'check-ins' table the rows that have as key the friend's ids. Moreover, there are certain columns that are selected from each row. These are the ones that are between two certain timestamps, as the columns have as qualifiers the timestamp of the check-in. The start timestamp is the one that corresponds to the UTC timestamp conversion of the start of the specific day at 12 am and the end timestamp corresponds to the end of the day at 12 am of the next day. In this way, the coprocessor returns only the check-ins of every friend that were made the intended day.
5. The client merges the the results that he got back from the multiple coprocessor calls into one list of user check-ins. Finally, the client sorts this list in order to present the check-ins of the user's friends into chronological order.

5.1.3 Correlated Most Visited POI query

This query contains the following question to the available data:

”Get the number of times that a user's friends have visited the user's most visited POI”

The query is implemented in the following steps:

1. The client calls the coprocessor that returns the friends of the specified user. The coprocessor runs on the region server that contains the row of the table 'friends' that has as key the user id of the desired user. The client receives back a list of user id's that represent his friends.
2. The client splits the list of friends into sections according to the initial split of the 'check-ins' table into 32 regions. In this way, each splitted friends list will contain row keys that belong to only one region server.
3. The client calls the most visited POI coprocessor in order to get the most visited POI of the specific user.
4. The client issues a call to the coprocessor that calculates the correlated most visited POI for every splitted friend list. More specifically, the client starts a new thread that is responsible for calling the coprocessor and getting the result back. In this way, the client issues parallel calls to the coprocessor and the calculations to the respective region servers are done simultaneously.

5. The coprocessor that runs on each region server gets from the 'check-ins' table the rows with the check-ins of all the friends that are on that region. Then he iterates over every check-in and checks whether it's location is the same with the location of the most visited POI of the original user. In this way it counts how many times each friend went to that place and returns that counter as a result.
6. The client merges the results that he got back from the multiple coprocessor calls and presents to the user which friends of him went to his most visited POI and how many times they visited that place.

5.2 Queries execution

5.2.1 Workload

Using the above queries we created a workload in order to test the behavior of the HBase cluster to multiple requests. More specifically, the workload consists of the three different type of queries, since they all refer to the same HBase tables. Moreover, the workload takes as input the number of queries that will be executed. In addition, since all queries include the retrieval of the friends of one user, that user is chosen randomly. The user that corresponds to each query is chosen using a generator of random numbers that follows a uniform distribution. The range of this distribution is the total number of available users, which is 9464.

Then, according to the number of queries, all types of queries participate in the workload in a cyclic assignment. For example, if the client wants 5 queries to be executed then those will be 1. most visited POI, 2. news feed, 3. corellated most visited POI, 4. most visited POI, 5. news feed. Also, the queries are executed in parallel as different threads. In this way, the HBase cluster receives simultaneously the query requests. The total execution time of the workload will be the biggest execution time amongst the queries of the workload.

5.2.2 Workload generation setup

There is one client that according to the specified number of queries to be executed, creates the workload in the way described previously. This client is hosted on a VM with 2 CPU, 4 GB RAM and 40 GB disk. The client is responsible to receive the number of queries to be executed and create the workload in the way described previously. The queries arrive at the same time in the HBase cluster that was described in section 4.1.

The workload is executed in the following setup:

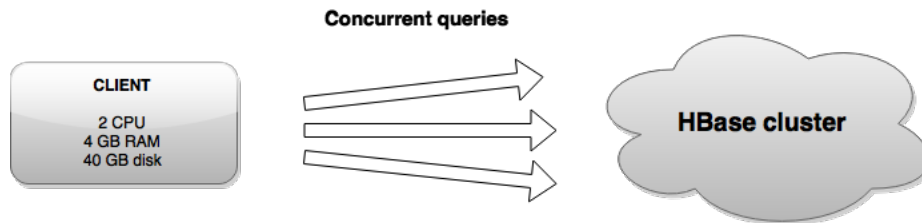


Figure 5.1: Query execution architecture

5.2.3 Scalability testing

Using the above architecture we performed a scalability testing, in order to evaluate how HBase handles the above workload for the specified dataset storage model in different cluster sizes. We calculated the latency and throughput of the system. More specifically, latency is the mean execution time of the queries. Throughput is the number of queries executed per second.

$$latency = \frac{\text{total execution time of all queries}}{\text{number of queries}} sec$$

$$throughput = \frac{\text{number of queries}}{\text{maximum query execution time}} queries/sec$$

We calculated the latency and throughput of the system for an increasing number of concurrent queries. We started with the HBase cluster having 32 nodes, as described previously. Then, we resized the cluster to 24, 16, 8 and 4 nodes in order to observe the variations in the latency and throughput. The restructure of the cluster was achieved by decommissioning the datanodes and region servers to the desired number. Both HDFS and HBase offer commands in order to achieve the resize of the cluster by moving data and regions into the remaining nodes, preventing data loss and ensuring that the regions will be data balanced.

The scalability testing can be presented in the following tables and corresponding plots:

Queries no.	32 nodes	24 nodes	16 nodes	8 nodes	4 nodes
5	2.4	2.8	2.4	2.6	3
10	2.5	2.7	2.7	3.2	4.6
15	3.2	3.06	3.6	3.73	5
20	3.4	3.9	4.2	4	5.95
25	5	5.32	5.56	5.76	6.56
30	5.5	5.76	6.36	6.03	7.66
35	6.37	6.08	7.06	7.51	9.91
40	6.97	7.52	8.05	8.35	10.725
45	7.55	8.31	9.02	10	10.86
50	8.34	9.02	10.32	11.2	13.6
55	8.96	10.14	11.21	11.45	13.85
60	10.26	10.88	11.3	12.58	15.98

Table 5.1: Scalability testing - Latency

The latency increases as the number of nodes of the HBase cluster decreases. Latency is the mean execution time of the queries. Thus, it is expected that the latency will increase when there are fewer nodes in the cluster. In more details, when the cluster size reduces there are fewer servers to handle the read requests and calculations that accompany the query. Therefore, the remaining HBase region servers have to handle more workload as opposed to when the cluster size was bigger. More specifically, since the data on the tables 'friends' and 'check-ins' are splitted into 32 regions in order to be equally distributed in the 32 nodes cluster, when we reduce the nodes, a region server can end up with more than one region from each table. In this way, the region server will have to handle more coprocessor calls as he has to server calculations for more data than before. Consequently, the execution time of the queries increases as there are fewer region servers and datanodes. Additionally, the latency increases as the client sends more concurrent queries. This happens due to the fact that as the number of concurrent queries elevate the servers cannot resolve them simultaneously and many of them have to wait in the servers queues. Therefore, there are queries that have additional waiting time to their total execution time. This waiting time becomes bigger when there are fewer servers to handle the concurrent queries. On the same level, when there is a small number of concurrent queries on different cluster sizes, then the latency is approximately the same due to the fact that there is no added waiting time.

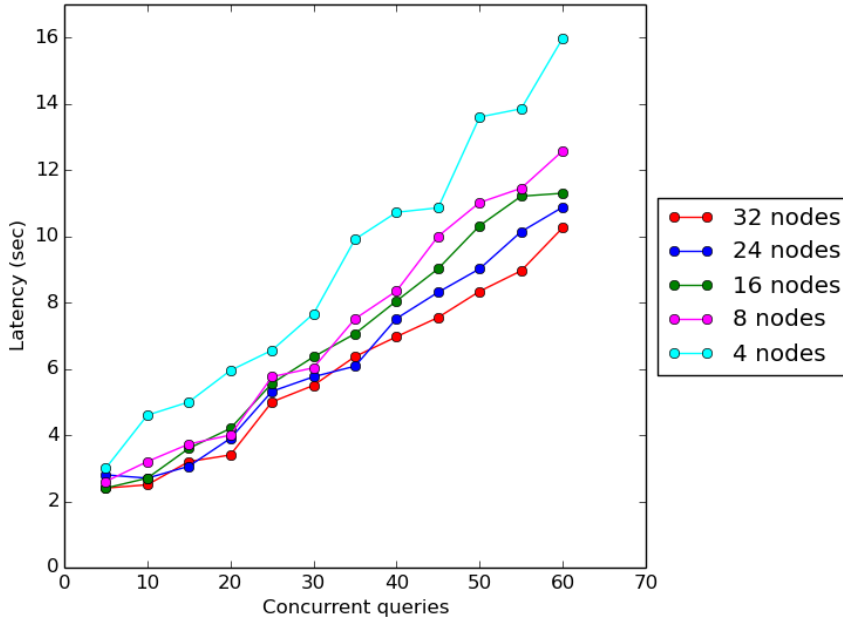


Figure 5.2: Scalability testing - Latency

Queries no.	32 nodes	24 nodes	16 nodes	8 nodes	4 nodes
5	1.66	1.66	1.66	1.66	1.25
10	3.33	2.5	2	1.8	1.66
15	3.75	3.45	3.2	2.9	2.54
20	4	3.9	3.5	3	2.5
25	4.46	3.77	3.46	3.47	2.77
30	5	4.28	3.75	3.33	2.72
35	5	4.375	3.9	3.78	2.6
40	5.2	4.44	4	3.63	2.85
45	5.425	4.5	4.09	3.75	2.8125
50	5.55	4.54	4.16	3.84	2.63
55	5.5	4.58	4.13	3.66	2.65
60	5.2	4.61	4	3.52	2.72

Table 5.2: Scalability testing - Throughput

Throughput is the number of queries executed per second. It is indicative of the performance of the system. For the current setup and workload, we observe that the throughput increases as the size of the cluster becomes bigger. This is expected due to the fact that more servers can serve more concurrent queries. When the number of servers increases the amount of work assigned to each server reduces, if the cluster is balanced. In our case, we ensure that the cluster is balanced as far as the data are concerned, when we pre-split the data into 32 regions. Also, HBase runs a balancer that keeps the regions equally distributed to the region servers. Therefore, when the cluster has 32 region servers, those have to run calculations over fewer data when the coprocessors are called, as opposed to when the cluster has fewer nodes. Moreover, we can see that for a specific cluster size the throughput increases as the number of concurrent queries elevate. This indicates that the system can handle a bigger load without dropping in performance. However, there is a certain limit to the throughput that prevents it from increasing even more. In every real time system exists that limit and indicates that the system

cannot handle more queries per second for the specific cluster size. This limit is different according to the number of nodes, while a bigger cluster size can have bigger throughput. Finally, we can see that this limit in the throughput is not a fixed value for every cluster size. This is due to the fact that the workload is not preset and can vary according to the randomly selected user factor. Additionally, cache affects tremendously the performance. When data for a query are retrieved from the server's cache, then the response time is noticeably less than getting the data from the datanode's disk. Also, the workload is executed on a real time system, on a cluster of virtual machines whose performance can be affected by the rest users of the cloud service that accommodates the VMs. However, those variations to the maximum value of the throughput are small and don't prevent us from seeing a stable limitation to the system's throughput.

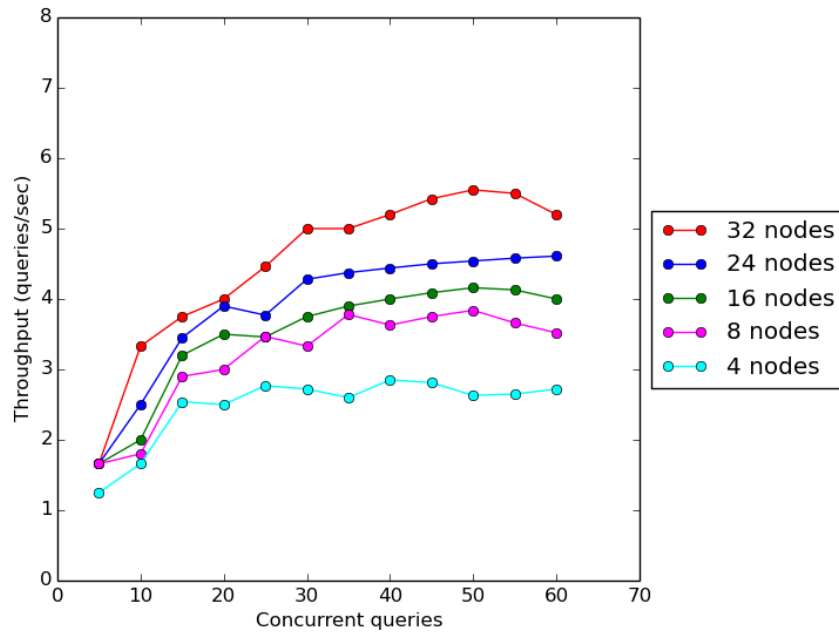


Figure 5.3: Scalability testing - Throughput

In conclusion, the HBase cluster that contains the social networking data produced by the generator and stored with the data model described previously, is scalable as it has the expected behavior of common scalable systems. The construction and execution of the workload that contains the three types of queries, allows as to test the scalability of the HBase cluster by calculating the latency and throughput.

Chapter 6

Conclusion

6.1 Summary

In this diploma thesis we were able to create a generator of spatio-temporal and textual social media data. We used real points of interest and ratings and reviews for these points as those were crawled by real users of the online travel service TripAdvisor. The utilization of these real points and the directions obtained from Google Directions API resulted in the creation of realistic routes, check-ins and GPS traces. Also, the fact that the generator uses random factors in order to take decisions about the place, duration and number of daily visits enhances the realistic aspect of the generated data. Additionally, most of these factors are instantiated by the corresponding input parameters, so that the generator can produce datasets with characteristics that vary. Therefore, we were able to create users that resemble real users found in well-known online social networking services.

The daily execution of the generator for a significant period of time resulted in the assemblance of a Big Data dataset of users with daily routes and check-ins to points of interest. More specifically, the generator created 9464 users, 1586537 check-ins and 38800019 GPS traces, which sum up to 3 GB data. In this way, we were able to re-create realistic social media data that we couldn't have access to due to privacy restrictions imposed by current social networking services.

With these data available, we were able to evaluate the distributed storage system that HBase provides and is used by a significant number of such social networking services. For this reason, we stored the created dataset into an HBase cluster using an appropriate data model and by ensuring that the cluster will be data balanced. Then, we implemented several queries over the available data that are common to such services. The queries were implemented using HBase coprocessors, which move the computation of intermediate results and other complex calculations to the region servers that contain the respective data, decongesting the client from a heavy computational load. We, then, created a workload with these types of queries in order to test the scalability of the HBase cluster. We measured the latency and throughput of the HBase cluster while changing the number of concurrent queries and cluster size. We came to the conclusion that the HBase cluster which contains the social networking data produced by the generator and stored with the specific data model, is scalable as it has the expected behavior of common scalable systems. Therefore, HBase provides indeed good performance and storage for Big Data social media services.

6.2 Related Work

There are already generators that produce social media data. One such generator is DATAGEN, which is a social network data generator, provided by LDBC-SNB (Linked Data Benchmark Council - Social Network Benchmark [11]). The LDBC-SNB aims to provide a realistic workload, consisting of a social network and a set of queries. DATAGEN is a synthetic data generator responsible for generating datasets for the three LDBC-SNB workloads: the Interactive, the Business Intelligence and the Analytical. DATAGEN mimics features of those found in real social network. They create a synthetic dataset due to the fact that it is difficult to find data with all the scaling characteristics their benchmark requires and collecting real data can be expensive or simply not possible due to privacy

concerns.

Also, there is LFR-Benchmark generator [18] which provides the network structure of a social network, being able to mimic some of the properties of real social media. The LFR generator produces social network graphs with power-law degree distribution and with implanted communities within the network.

6.3 Future Work

There are several ways in which the current diploma thesis could be extended. As far as the generator is concerned, the produced data could be enhanced by also adding more textual data such as Facebook posts or Twitter tweets. Also, the generator could use a photo and video pool in order to simulate the upload of photos and videos, which is a basic functionality of most social networking services. As far as the generator execution is concerned, the generator could be instantiated with other values in order to create a more polymorphic dataset. Also, the usage of Google Maps API for Work enables 100000 directions requests per 24 hour period, which is significantly more than the limit of 2500 daily requests for the users of the free API. Therefore, the generator can be used more efficiently for users subscribed to the Google Work services.

Moreover, regarding the implemented queries, there are many more that can be implemented. Also, there could be queries over the data of GPS traces which are created by the generator. In this way, the created workload can be more efficient into the scalability testing of the HBase cluster. Finally, other distributed storage and management systems could be evaluated in a similar way, such as Hadoop.

Bibliography

- [1] Twitter statistics. <https://about.twitter.com/company>. [Online, accessed June 2015].
- [2] Google directions api. <https://developers.google.com/maps/documentation/directions/>. [Online, accessed June 2015].
- [3] Encoded polyline algorithm format. <https://developers.google.com/maps/documentation/utilities/polylinealgorithm>. [Online, accessed June 2015].
- [4] Google maps static maps api. <https://developers.google.com/maps/documentation/staticmaps/>. [Online, accessed June 2015].
- [5] Twitter re-architect. <https://blog.twitter.com/2013/new-tweets-per-second-record-and-how>. [Online, accessed June 2015].
- [6] Apache hbase coprocessors. <https://hbase.apache.org/book.html#cp>. [Online, accessed June 2015].
- [7] Institutions using hbase. <http://wiki.apache.org/hadoop/Hbase/PoweredBy>. [Online, accessed June 2015].
- [8] Institutions using hadoop. <https://wiki.apache.org/hadoop/PoweredBy>. [Online, accessed June 2015].
- [9] Databases - wikipedia. <https://en.wikipedia.org/wiki/Database>. [Online, accessed June 2015].
- [10] R. Bayer and E. McCreight. Organization and maintenance of large ordered indices. In *Proceedings of the 1970 ACM SIGFIDET (Now SIGMOD) Workshop on Data Description, Access and Control*, SIGFIDET '70, pages 107–141, New York, NY, USA, 1970. ACM.
- [11] Peter Boncz. Ldbc: Benchmarks for graph and rdf data management. In *Proceedings of the 17th International Database Engineering & Applications Symposium*, IDEAS '13, pages 1–2, New York, NY, USA, 2013. ACM.
- [12] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. Bigtable: A distributed storage system for structured data. *ACM Trans. Comput. Syst.*, 26(2):4:1–4:26, June 2008.
- [13] E. F. Codd. A relational model of data for large shared data banks. *Commun. ACM*, 13(6):377–387, June 1970.
- [14] Antonin Guttman. R-trees: A dynamic index structure for spatial searching. In *Proceedings of the 1984 ACM SIGMOD International Conference on Management of Data*, SIGMOD '84, pages 47–57, New York, NY, USA, 1984. ACM.
- [15] Theo Haerder and Andreas Reuter. Principles of transaction-oriented database recovery. *ACM Comput. Surv.*, 15(4):287–317, December 1983.

- [16] Joseph M. Hellerstein, Jeffrey F. Naughton, and Avi Pfeffer. Generalized search trees for database systems. In *Proceedings of the 21th International Conference on Very Large Data Bases*, VLDB '95, pages 562–573, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
- [17] Ioannis Ioannou. Ανάλυση Συναισθήματος Σε Μεγάλο Όγκο Δεδομένων Κειμένου Με Χρήση Κατανεμημένων Τεχνικών Μηχανικής Εκμάθησης, July 2014. Diploma Thesis, National Technical University of Athens.
- [18] Andrea Lancichinetti, Santo Fortunato, and Filippo Radicchi. Benchmark graphs for testing community detection algorithms. *Physical Review E (Statistical, Nonlinear, and Soft Matter Physics)*, 78(4), 2008.
- [19] Ioannis Mytilinis, Ioannis Giannakopoulos, Ioannis Konstantinou, Katerina Doka, Dimitrios Tsitsigkos, Manolis Terrovitis, Lampros Giampouras, and Nectarios Koziris. Modissense: A distributed spatio-temporal and textual processing platform for social networking services. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, SIGMOD '15, pages 895–900, New York, NY, USA, 2015. ACM.
- [20] Regina Obe and Leo Hsu. *PostGIS in Action*. Manning Publications Co., Greenwich, CT, USA, 2011.
- [21] Regina Obe and Leo Hsu. *PostgreSQL: Up and Running*. O'Reilly Media, Inc., 2012.
- [22] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The hadoop distributed file system. In *Proceedings of the 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST)*, MSST '10, pages 1–10, Washington, DC, USA, 2010. IEEE Computer Society.