

UNIVERSITATEA BABEȘ-BOLYAI  
FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ

# **Modele de date NoSQL: contribuții la reprezentarea structurii și evaluarea performanței**

**Rezumatul tezei de doctorat**

Student-doctorand: Camelia-Florina ANDOR  
Coordonator științific: prof. dr. Bazil PÂRV

2021

# 1. Introducere

În contextul actual, alegerea unui sistem de gestiune a bazelor de date (SGBD) nepotrivit specificului unei aplicații are un impact semnificativ asupra performanței. Aplicațiile moderne generează o cantitate impresionantă de date, care necesită SGBD-uri adaptabile și performante, capabile să funcționeze cu ușurință într-un mediu distribuit. Aceste date generate într-un ritm accelerat se află în diverse formate și necesită procesări complexe. SGBD-urile relaționale nu au fost proiectate pentru a funcționa optim în asemenea condiții. Acest lucru a dus la apariția unei noi generații de SGBD-uri, adaptate nevoilor actuale de stocare și de procesare a datelor. Aceste noi SGBD-uri fac parte din categoria NoSQL, sunt în general proiecte de tipul *open source* și funcționează cu succes într-un mediu distribuit. Există în prezent o multitudine de implementări NoSQL disponibile pe piață, fiecare având un anumit model de date, un limbaj de interogare propriu și funcționalități automatizate de replicare și distribuție a datelor. Interacțiunea cu aceste SGBD-uri este mai dificilă, acest fapt datorându-se particularităților fiecărei implementări dar și naturii puternic specializate în rezolvarea anumitor probleme din sfera *big data*. Regulile de modelare a datelor diferă de la o implementare la alta, iar familiarizarea cu acestea necesită destul de mult timp și efort. Unele implementări NoSQL nu beneficiază de o interfață grafică de tip client, iar interacțiunea cu acestea se face doar prin intermediul unei aplicații de tip consolă. Dezvoltatorii SGBD-urilor NoSQL sunt în general foarte preocupați de adăugarea unor noi funcționalități și de îmbunătățirea celor existente, astfel că acestea au o evoluție rapidă și prezintă diferențe semnificative între versiuni. Gradul înalt de performanță și flexibilitatea structurii datelor sunt două dintre cele mai mari avantaje ale SGBD-urilor NoSQL.

Alegerea SGBD-ului NoSQL potrivit pentru o aplicație dată este o sarcină destul de greu de realizat. Marea varietate de produse disponibile și diferențele foarte mari dintre ele îngreunează semnificativ realizarea unei comparații corecte. Documentațiile tehnice ale produselor sunt de asemenea foarte complexe și necesită un studiu îndelungat. Experimentele de evaluare a performanței SGBD-urilor NoSQL oferă informații importante despre performanța acestora în contexte specifice, fiind utile în procesul de alegere a SGBD-ului potrivit.

În această lucrare ne-am propus să simplificăm procesul de alegere a SGBD-ului NoSQL potrivit unei aplicații sau unui tip de aplicații date și să facilităm interacțiunea utilizatorilor cu bazele de date NoSQL. Specialiști din toate domeniile lucrează în prezent cu cantități mari de date care necesită procesări complexe. Aplicațiile care facilitează interacțiunea cu baza de date sunt de mare ajutor, în special în cazul utilizatorilor care nu sunt specializați în domeniul informaticii.

Contextul și descrierea problemelor pe care ne-am propus să le soluționăm în lucrarea de față sunt prezentate în capitolele 1, 2 și 3.

Contribuțiile originale ale lucrării sunt prezentate în capitolele 4, 5 și 6.

## 2. Structura tezei de doctorat

Teza de doctorat este compusă din șapte capitole.

În capitolul 1 sunt prezentate neajunsurile modelului de date relațional, care au condus la apariția SGBD-urilor NoSQL, dar și cele patru modele de date NoSQL fundamentale. SGBD-urile NoSQL folosesc modele de date diferite de cel relațional, sunt proiectate pentru a funcționa fără probleme într-un mediu distribuit și oferă un grad mai mare de flexibilitate a structurii datelor. Fiecare model de date NoSQL fundamental (*key-value*, *document*, *column-family*, *graph*) este descris într-o secțiune separată. De asemenea, sunt abordate și alte teme importante din zona SGBD-urilor NoSQL, cum ar fi replicarea și distribuția datelor.

Capitolul 2 abordează problema dificultății alegerii SGBD-ului NoSQL potrivit pentru o aplicație dată. Diferențele mari dintre SGBD-urile NoSQL, datorate modelelor de date diferite, limbajelor de interogare proprii și altor caracteristici îngreunează realizarea unei comparații corecte între diferite implementări NoSQL. Experimentele de evaluare a performanței SGBD-urilor NoSQL reprezintă o soluție în acest caz, deoarece dezvăluie informații importante legate de performanța fiecărei implementări NoSQL testate în condiții specifice. Luând în considerare atât cazurile de utilizare comune, cât și gradul de complexitate al modelului de date, sunt alese două modele de date NoSQL considerate mai versatile (*document* și *column-family*) și este descrisă câte o implementare NoSQL reprezentativă pentru fiecare (MongoDB, respectiv Apache Cassandra). În continuare, este abordat subiectul aplicațiilor de evaluare a performanței. Sunt descrise atât tipurile de aplicații de evaluare a performanței cât și rolul lor într-un experiment de evaluare a performanței. Cea mai populară aplicație de evaluare a performanței, *YCSB*, este de asemenea prezentată. Sunt descrise pe scurt și studii experimentale de evaluare a performanței din literatură în care a fost folosită aplicația *YCSB*.

În capitolul 3 este abordată problema interacțiunii dificile cu baze de date NoSQL, datorată în principal faptului că fiecare SGBD NoSQL are un limbaj de interogare propriu. O soluție la această problemă ar putea fi aplicațiile de tip *Visual Query System*, care folosesc o reprezentare conceptuală a structurii bazei de date și a interogărilor și permit utilizatorilor interacțiunea cu baza de date într-o manieră vizuală. Sunt prezentate diferite abordări din literatura de specialitate care folosesc grafuri conceptuale pentru reprezentarea structurii bazei de date și a interogărilor. Aceste abordări vizează bazele de date relaționale și datele stocate în format XML și sunt implementate în aplicații de tip *Visual Query System*.

Capitolul 4 prezintă cinci studii de caz bazate pe experimente de evaluare a performanței SGBD-urilor MongoDB și Apache Cassandra. Metricile de performanță analizate în aceste studii sunt *throughput*, *average latency* (latența medie) și *total runtime* (timpul total de execuție). În ultimele două studii de caz sunt folosite versiuni noi ale SGBD-urilor MongoDB și Apache Cassandra. În urma analizei rezultatelor ultimelor două studii de caz, s-a dovedit că există diferențe semnifica-

tive ale performanței între versiuni diferite ale aceluiași SGBD. Aplicația de evaluare a performanței folosită în acele cinci studii de caz este *YCSB*. În urma utilizării aplicației *YCSB*, am remarcat anumite neajunsuri ale acesteia, precum lipsa unei interfețe grafice bazate pe un model conceptual de reprezentare a structurii bazei de date și a operațiilor implicate în teste. Pentru realizarea unor experimente de evaluare a performanței care necesită seturi de date și seturi de operații personalizate cu *YCSB* sunt necesare modificări sau extinderi ale codului sursă.

În capitolul 5 propunem două metode de reprezentare a structurii bazei de date și a interogărilor bazate pe grafuri conceptuale pentru MongoDB, respectiv Apache Cassandra. Aceste metode pot fi implementate atât în aplicații de tip *Visual Query System*, cât și într-o aplicație de evaluare a performanței SGBD-urilor.

Capitolul 6 prezintă două implementări ale metodei de reprezentare a structurii bazei de date și a interogărilor MongoDB. Aplicația web *Conceptual Graphs for MongoDB* se încadrează în categoria *Visual Query System* și oferă posibilitatea de a reprezenta structura bazei de date și a interogărilor într-o manieră vizuală cu ajutorul grafurilor conceptuale. Aplicația transformă reprezentările conceptuale în instrucțiunile echivalente din MongoDB, care pot fi mai apoi executate. Aplicația desktop *DBMark* este o aplicație de evaluare a performanței care oferă o interfață grafică bazată pe grafuri conceptuale. Momentan, aplicația *DBMark* poate fi folosită pentru evaluarea performanței SGBD-ului MongoDB. Cu ajutorul acestei aplicații se pot realiza experimente de evaluare a performanței în care se folosesc seturi de date și seturi de operații personalizate. Interfața grafică bazată pe grafuri conceptuale facilitează personalizarea experimentelor, iar aplicația poate genera automat reprezentarea conceptuală a structurii unei baze de date deja existente. Bineînțeles, se pot crea în aplicație reprezentări conceptuale noi și în acest caz se vor genera în mod automat instrucțiunile MongoDB echivalente.

În ultimul capitol sunt prezentate concluziile acestei lucrări și direcțiile de cercetare viitoare. Ne propunem ca în viitor să dezvoltăm metode de reprezentare a structurii bazei de date și a interogărilor bazate pe grafuri conceptuale pentru alte SGBD-uri NoSQL, dar și să realizăm alte studii de caz bazate pe experimente de evaluare a performanței SGBD-urilor NoSQL.

### 3. Contribuții originale

Cele două probleme abordate în această lucrare sunt reprezentate de dificultatea și complexitatea procesului de alegere a SGBD-ului NoSQL potrivit pentru o aplicație dată și dificultatea interacțiunii cu SGBD-urile NoSQL, care se datorează în principal diferențelor semnificative dintre limbajele de interogare. Contribuțiile originale vizează simplificarea procesului de alegere a SGBD-ului NoSQL potrivit unei aplicații date și facilitarea interacțiunii cu SGBD-uri NoSQL în cazul utilizatorilor care nu sunt specializați în informatică.

Experimentele de evaluare a performanței au un rol important în simplificarea procesului de alegere a SGBD-ului NoSQL potrivit pentru o aplicație dată. Acestea oferă informații importante despre performanța SGBD-urilor evaluate în condițiile specifice cazului de utilizare considerat. În capitolul 4 sunt prezentate cinci studii de caz realizate pe baza experimentelor de evaluare a performanței. În aceste studii de caz a fost analizată performanța a două dintre cele mai populare SGBD-uri NoSQL, MongoDB și Apache Cassandra. Metricile de performanță analizate sunt *throughput*, *average latency* (latența medie) și *total runtime* (timpul total de execuție). Primele trei studii de caz (SC1, SC2 și SC3) analizează performanța SGBD-urilor MongoDB și Apache Cassandra din perspectiva *throughput*-ului, a latenței medii și a timpului total de execuție. A fost realizată și analiza statistică a rezultatelor obținute. SGBD-urile NoSQL au o evoluție rapidă, și există diferențe semnificative între versiuni diferite ale aceluiași SGBD care pot avea impact asupra performanței. În ultimele două studii de caz (SC4 și SC5) au fost refăcute experimentele din SC1 și SC3 folosind versiuni noi ale SGBD-urilor, iar în urma analizei rezultatelor s-a dovedit că impactul asupra performanței este semnificativ.

Experimentele de evaluare a performanței SGBD-urilor NoSQL sunt realizate cu ajutorul unor aplicații de evaluare a performanței, numite și *benchmarking tools* sau *benchmarking frameworks*. În capitolul 2 au fost prezentate cele două tipuri de aplicații de evaluare a performanței și rolul pe care îl au acestea în procesul de evaluare a performanței. În realizarea experimentelor care stau la baza celor cinci studii de caz prezentate în capitolul 4 a fost folosită cea mai populară aplicație de evaluare a performanței SGBD-urilor NoSQL, *YCSB*. O descriere a acestei aplicații dar și a SGBD-urilor MongoDB și Apache Cassandra a fost realizată în capitolul 2. În urma utilizării aplicației *YCSB* în toate cele cinci studii de caz, am identificat anumite neajunsuri ale acesteia cu privire la personalizarea experimentelor de evaluare a performanței, cum ar fi lipsa unei interfețe grafice care să permită o reprezentare la nivel conceptual a structurii bazei de date și a operațiilor executate într-un test de performanță.

În capitolul 5 am propus două metode de reprezentare a structurii bazei de date și a interogărilor cu ajutorul grafurilor conceptuale pentru SGBD-urile MongoDB și Apache Cassandra. Acestea pot fi implementate într-o aplicație de evaluare a performanței SGBD-urilor MongoDB și Apache Cassandra care simplifică efectuarea experimentelor de evaluare a performanței ce folosesc seturi de date

și de operații personalizate. Metoda de reprezentare a structurii bazei de date și a interogărilor MongoDB a fost implementată în aplicația *DBMark*, care este descrisă în capitolul 6. Aplicația *DBMark* oferă o interfață grafică care permite reprezentarea structurii bazei de date și a interogărilor MongoDB cu ajutorul grafurilor conceptuale și poate fi folosită în efectuarea experimentelor de evaluare a performanței care implică seturi de date și de operații personalizate.

Interacțiunea cu SGBD-urile NoSQL, care este îngreunată de marea varietate de limbaje de interogare, poate fi simplificată cu ajutorul aplicațiilor de tip *Visual Query System*. Acestea facilitează interacțiunea cu baza de date, cu precădere în cazul utilizatorilor care nu sunt specialiști în domeniul informaticii. Într-o aplicație de tipul *Visual Query System*, utilizatorul poate reprezenta structura bazei de date și a interogărilor într-o manieră vizuală, folosind un limbaj de modelare conceptual. Metodele de reprezentare a structurii bazei de date și a interogărilor pentru SGBD-urile MongoDB și Apache Cassandra propuse în capitolul 5 pot fi implementate și în aplicații de tip *Visual Query System*. Metoda de reprezentare a structurii bazei de date și a interogărilor MongoDB a fost implementată în aplicația *Conceptual Graphs for MongoDB*, descrisă în capitolul 6. Aplicația *Conceptual Graphs for MongoDB* simplifică interacțiunea cu baza de date în cazul SGBD-ului MongoDB, transformând reprezentările bazate pe grafuri conceptuale a structurii bazei de date și a interogărilor în instrucțiuni ce pot fi executate în MongoDB.

În ceea ce privește direcțiile viitoare de cercetare, ne propunem să implementăm metoda de reprezentare a structurii bazei de date și a interogărilor Apache Cassandra într-o aplicație de tipul *Visual Query System*, dar și în aplicația de evaluare a performanței *DBMark*, pentru a putea fi ulterior folosită și în evaluarea performanței SGBD-ului Apache Cassandra. În ceea ce privește experimentele de evaluare a performanței SGBD-urilor NoSQL, ne propunem să efectuăm noi experimente care să includă și alte SGBD-uri NoSQL și să utilizăm o infrastructură *cloud-based*. Un alt obiectiv viitor este dezvoltarea unor metode de reprezentare a structurii bazei de date și a interogărilor cu ajutorul grafurilor conceptuale pentru alte SGBD-uri NoSQL. Acestea ar putea fi ulterior folosite pentru a extinde funcționalitățile aplicației de evaluare a performanței *DBMark* și pentru a implementa noi aplicații din categoria *Visual Query System*.

**Cuvinte-cheie:** baze de date NoSQL; aplicații de evaluare a performanței; MongoDB; Apache Cassandra; YCSB; grafuri conceptuale; modele de date; big data.

# Cuprinsul tezei de doctorat

<b>Lista figurilor</b>	<b>5</b>
<b>Lista tabelelor</b>	<b>5</b>
<b>Lista publicațiilor</b>	<b>6</b>
<b>Introducere</b>	<b>8</b>
<b>1 Modele de date NoSQL</b>	<b>11</b>
1.1 Modelul de date relațional . . . . .	11
1.2 Modele de date NoSQL . . . . .	12
1.2.1 Modelul Key-Value . . . . .	13
1.2.2 Modelul Document . . . . .	14
1.2.3 Modelul Column-family . . . . .	15
1.2.4 Modelul Graph . . . . .	16
1.2.5 Replicarea și distribuția datelor . . . . .	17
1.3 Concluzii . . . . .	17
<b>2 Sisteme de gestiune a bazelor de date NoSQL și aplicații de evaluare a performanței</b>	<b>19</b>
2.1 Sisteme de gestiune a bazelor de date NoSQL . . . . .	19
2.1.1 MongoDB . . . . .	19
2.1.2 Apache Cassandra . . . . .	21
2.2 Aplicații de evaluare a performanței pentru sisteme de gestiune a bazelor de date NoSQL	22
2.3 Concluzii . . . . .	24
<b>3 Grafuri conceptuale</b>	<b>26</b>
3.1 Utilizarea grafurilor conceptuale în domeniul bazelor de date . . . . .	26
3.2 Exemple de grafuri conceptuale . . . . .	28
3.2.1 Content Management System . . . . .	28
3.2.2 Blogging System . . . . .	29
3.2.3 Aplicație de comerț electronic . . . . .	31
3.2.4 DBLP . . . . .	32
3.3 Concluzii . . . . .	34

<b>4</b>	<b>Evaluarea performanței sistemelor de gestiune a bazelor de date NoSQL</b>	<b>35</b>
4.1	SC1: Evaluarea throughput sub Windows 7 . . . . .	38
4.2	SC2: Evaluarea latenței medii sub Windows 7 . . . . .	43
4.3	SC3: Evaluarea timpului total de execuție sub Windows 7 . . . . .	46
4.4	SC4: Evaluarea throughput sub Windows 10 . . . . .	52
4.5	SC5: Evaluarea timpului total de execuție sub Windows 10 . . . . .	58
4.6	Concluzii . . . . .	65
<b>5</b>	<b>Reprezentarea structurii bazei de date și a interogărilor NoSQL cu ajutorul grafurilor conceptuale</b>	<b>67</b>
5.1	Reprezentarea structurii bazei de date și a interogărilor MongoDB cu ajutorul grafurilor conceptuale . . . . .	68
5.1.1	Reprezentarea grafică a structurii datelor în MongoDB cu ajutorul grafurilor conceptuale . . . . .	68
5.1.2	Exemplu de reprezentare a structurii bazei de date în MongoDB cu ajutorul grafurilor conceptuale . . . . .	69
5.1.3	Reprezentarea grafică a interogărilor MongoDB cu ajutorul grafurilor conceptuale . . . . .	71
5.2	Reprezentarea structurii bazei de date și a interogărilor Cassandra cu ajutorul grafurilor conceptuale . . . . .	75
5.2.1	Reprezentarea grafică a structurii datelor în Cassandra cu ajutorul grafurilor conceptuale . . . . .	76
5.2.2	Exemplu de reprezentare a structurii bazei de date în Cassandra cu ajutorul grafurilor conceptuale . . . . .	79
5.2.3	Reprezentarea grafică a interogărilor Cassandra cu ajutorul grafurilor conceptuale . . . . .	81
5.3	Concluzii . . . . .	85
<b>6</b>	<b>Aplicații</b>	<b>86</b>
6.1	Aplicație de tip Visual Query System pentru MongoDB . . . . .	86
6.1.1	Tehnologii folosite în dezvoltarea aplicației . . . . .	86
6.1.2	Analiza cerințelor . . . . .	86
6.1.3	Descrierea aplicației . . . . .	87
6.2	Aplicație de tip benchmarking framework pentru MongoDB . . . . .	89
6.2.1	Tehnologii folosite în dezvoltarea aplicației . . . . .	90
6.2.2	Analiza cerințelor . . . . .	90
6.2.3	Descrierea aplicației . . . . .	91
6.2.4	Un exemplu de utilizare a aplicației cu un set de date personalizat . . . . .	92
6.3	Concluzii . . . . .	103
<b>7</b>	<b>Concluzii și direcții viitoare de cercetare</b>	<b>105</b>
	<b>Bibliografie</b>	<b>108</b>