



POLITECHNIKA WARSZAWSKA

Wydział Elektroniki i Technik Informacyjnych

Instytut Informatyki

Rok akademicki 2013/2014

PRACA DYPLOMOWA MAGISTERSKA

Jakub Turek

Mechanizm modelowania danych i mapowania obiektowego dla Apache Cassandra

Praca wykonana pod kierunkiem
dra inż. Jakuba Koperwasa

Ocena:

.....

*Podpis Przewodniczącego Komisji
Egzaminu Dyplomowego*

Kierunek: Informatyka
Specjalność: Inżynieria Systemów Informatycznych
Data urodzenia: 1990.01.09
Data rozpoczęcia studiów: 2013.02.20

Życiorys

Urodziłem się 9 stycznia 1990 roku w Łodzi. W 1997 roku rozpocząłem edukację w Szkole Podstawowej nr 7 w Łodzi. W latach 2003-2006 kontynuowałem naukę w Gimnazjum nr 42 im. Władysława Stanisława Reymonta w Łodzi. Od 2006 roku uczyłem się w Liceum Ogólnokształcącym nr 31 im. Ludwika Zamenhofs w Łodzi. W 2009 roku zdałem egzaminy maturalne i ukończyłem szkołę licealną z wyróżnieniem. W latach 2009-2013 studiowałem dziennie informatykę na Wydziale Elektroniki i Technik Informacyjnych Politechniki Warszawskiej. Ukończyłem studia z wynikiem celującym i odebrałem tytuł zawodowy inżyniera. Obecnie kończę pracę dyplomową masterską pod kierownictwem Instytutu Informatyki. We wrześniu 2012 roku rozpocząłem pracę zawodową jako programista aplikacji do zarządzania procesami biznesowymi oraz aplikacji mobilnych w firmie Xentivo, gdzie pracuję do dziś. Moją pasją jest tworzenie aplikacji mobilnych oraz internetowych, które uruchamiane są w środowisku iOS.

.....
Podpis studenta

Egzamin dyplomowy:

Złożył egzamin dyplomowy w dniu:

z wynikiem:

Ogólny wynik studiów:

Dodatkowe uwagi i wnioski Komisji:

.....

Streszczenie

Celem pracy dyplomowej jest stworzenie mechanizmu do modelowania i mapowania obiektowego dla aplikacji wykorzystujących bazę danych Apache Cassandra. System powinien cechować się wysoką wydajnością i wsparciem dla wzorców modelowania. Dodatkowym celem jest zachowanie zgodności z istniejącymi mechanizmami mapowania dla baz relacyjnych. Nie wszystkie cele udało się zrealizować. Stworzenie wydajnego mapowania wymagało odejścia od relacyjnego modelu danych. Praca dyplomowa zawiera zbiór informacji teoretycznych, badań oraz studia przypadków, które posłużyły do nakreślenia pryncypiów systemu modelowania obiektowego. Obejmuje również implementację omawianych mechanizmów.

Data modeling and object mapping mechanism for Apache Cassandra

Summary

The goal of this thesis is to create a modeling and object mapping mechanism for applications which use Apache Cassandra database. The mechanism should provide good performance and be capable of making use of Cassandra design patterns. An additional goal is to explore possibility of using existing interfaces for object-relational mapping. Not all goals could be achieved. Creating a fast system required abandoning object-relational mapping compatibility. The thesis contains of theoretical information, research results and case studies which contributed to specify principles of the mechanism. It also provides functional implementation of described ideas.

Spis treści

1	Wstęp	1
1.1	Zakres pracy	1
1.2	Motywacja	2
1.3	Zawartość pracy	3
2	Apache Cassandra	5
2.1	Model danych	6
2.2	Dystrybucja danych	7
2.3	Algorytmy zapisu/usuwania danych	8
2.4	Obszary zastosowania Cassandra	10
2.5	Struktura danych a modelowanie	16
2.6	CQL	18
2.7	Przykład modelowania dziedziny danych	19
2.7.1	Twissandra	20
2.8	Modelowanie - wzorce i antywzorce	29
3	Mapowanie obiektowo-relacyjne	31
3.1	Java Persistence API	33
3.2	Kundera	35
3.2.1	Wydajność	35
3.3	Wnioski z badania wydajności	38
3.4	Hibernate OGM	40
3.5	Koncepcja mapowania dla Cassandra	41

4	Modelowanie obiektowe dla Cassandra	44
4.1	Object Modeling for Cassandra	46
4.2	OMC - podstawowe pojęcia	48
4.2.1	Model	48
4.2.2	Pole	48
4.2.3	Silnik	48
4.3	Definiowanie modelu	49
4.4	Modelowanie zależności między danymi	51
4.4.1	Denormalizacja przez pole	54
4.4.2	Denormalizacja przez tabelę	56
4.4.3	Normalizacja przez tabelę	58
4.4.4	Porównanie wydajności metod zarządzania zależnościami	59
4.4.5	Wnioski z analizy wydajności modelowania zależności .	64
4.5	Wsparcie dla wzorców modelowania	68
4.5.1	Szereg zdarzeń	69
4.5.2	Wsparcie dla kolejek	72
4.5.3	Selektywna aktualizacja	74
4.5.4	Indeks wartości unikalnych	76
4.6	Przetwarzanie partiami	78
4.7	Wsparcie dla liczników	80
4.8	Migracje pomiędzy wersjami modelu	81
4.8.1	Wydajność migracji danych	84
4.8.2	Potencjalne problemy migracji	87
4.9	Profilowanie i testowanie modelu danych	87
4.9.1	Zasilanie danymi	88
4.10	Lista życzeń użytkownika - studium przypadku	91
5	Podsumowanie	94
	Bibliografia	96

Rozdział 1

Wstęp

Tematem pracy jest stworzenie systemu mapowania obiektowego dla bazy danych Apache Cassandra. Celem pracy jest zweryfikowanie możliwości wykorzystania interfejsów analogicznych mechanizmów stworzonych dla relacyjnych baz danych. Podstawowym wymaganiem jest zachowanie wysokiej wydajności zapisu, która wyróżnia Cassandrę na tle innych systemów bazodanowych, a także możliwość czerpania z praktyk zwiększających efektywność modelowania opisanych przez użytkowników Cassandry. Niniejsza praca prezentuje zbiór rozważań teoretycznych, badań oraz doświadczeń z użytkowania, które przyczyniły się do powstania takiego mechanizmu. Częścią pracy jest również implementacja mechanizmu.

1.1 Zakres pracy

Zakres pracy obejmuje następujące elementy:

1. Badanie możliwości wykorzystania interfejsów mapowania obiektowo-relacyjnego dla bazy danych Cassandra:
 - Wyszukanie istniejących implementacji.
 - Pomiar wydajności ze szczególnym naciskiem na różnice w stosunku do relacyjnych baz danych oraz modelowania i komunikacji

z wykorzystaniem języka zapytań Apache Cassandra.

2. Zaproponowanie własnego interfejsu zorientowanego na wysoką wydajność oraz wsparcie w modelowaniu danych z wykorzystaniem najlepszych praktyk:

- Zebranie i opisanie sposobów modelowania zależności między danymi oraz porównanie ich wydajności w różnych przypadkach użycia.
- Zebranie i opisanie zbioru najlepszych wzorców modelowania oraz możliwości wspierania ich w mechanizmie.
- Wykonanie referencyjnej implementacji zaproponowanego interfejsu.
- Przeprowadzenie badań wydajnościowych implementacji.

1.2 Motywacja

W dobie szybkiego postępu technologicznego i intensywnego rozwoju informatyki, a zwłaszcza powszechnego dostępu do sieci Internet oraz niskich kosztów składowania danych, zaczął rozwijać się trend nazywany przez specjalistów od marketingu *big data*¹. *Big data* to termin używany do określania dużych rozmiarów, różnorodnych i często zmieniających się zbiorów danych. [1] Efektywne przetwarzanie takich danych wymaga stosowania innowacyjnych, stale ulepszanych rozwiązań technologicznych.

Powstanie ogromnych portali społecznościowych obsługiwanych przez dziesiątki tysięcy fizycznych maszyn spowodowało pojawienie się wyzwań, które wcześniej nie były brane pod uwagę. Awarie maszyn, dotychczas traktowane jako sytuacje wyjątkowe, przy tej skali użytkownia stały się regułą. W połączeniu z wymaganiami na niskie czasy obsługi setek tysięcy jednoczesnych żądań doprowadziło to do osiągnięcia kresu możliwości wykorzystywanych od

¹Z angielskiego oznacza to dosłownie „wielkie dane”.

wielu lat relacyjnych baz danych. Aby móc sprostać tym warunkom zaczęły powstawać nowe silniki bazodanowe, które odchodziły od postulatów transakcyjności i klasycznej, tabelarycznej reprezentacji danych. Jednym z najlepszych rozwiązań do obsługi wysoce rozproszonych środowisk jest Apache Cassandra.

Negatywną stroną wprowadzenia nowych rozwiązań z zakresu przechowywania danych było odcięcie się od wielu mechanizmów, które bardzo ułatwiały pracę z aplikacjami. Przede wszystkim nie było możliwe użycie istniejących systemów mapowania obiektowo-relacyjnego. Ich wykorzystanie stało się na tyle powszechne, że w dniu dzisiejszym mogą być uznawane za standard w tworzeniu aplikacji bazodanowych.

Wraz z rozwojem Cassandry, a zwłaszcza wprowadzeniem języka zapytań o składni przypominającej ten z relacyjnych baz danych, pojawiła się możliwość ponownego wykorzystania istniejących interfejsów. Nawet jeżeli zachowanie wysokiej wydajności wymaga wprowadzenia pewnych modyfikacji, to warto to uczynić. Dzięki temu możliwe jest zapewnienie narzędzia do szybszego rozwoju aplikacji oraz zmniejszenie bariery wejścia związanej z opanowaniem podstaw nowego systemu.

Autor ma nadzieję, że przeprowadzone przez niego badania przyczynią się do rozwoju prac nad efektywnym modelowaniem danych w Apache Cassandra oraz pomogą wypełnić lukę w oprogramowaniu wspomagającym pracę z nierelacyjnymi bazami danych. Ponadto Autor wyraża opinię, że rozwój dedykowanych systemów mapowania przyczyni się do szerszego wykorzystywania przedmiotowego systemu bazodanowego.

1.3 Zawartość pracy

Rozdział 2 opisuje podstawy teoretyczne związane z silnikiem Apache Cassandra, które są niezbędne do zrozumienia całej pracy. W sekcji 2.1 Autor przedstawia stosowany w silniku model danych. Sekcja 2.2 opisuje sposób dystrybucji wierszy w klastrze. W punkcie 2.5 zostały omówione, na przykła-

dzie, różnice pomiędzy relacyjnym modelem danych a schematem Cassandra. Sekcja 2.6 opowiada o języku zapytań CQL, natomiast w punkcie 2.8 zostało wprowadzone pojęcie wzorców i antywzorców modelowania.

W rozdziale 3 zostały przedstawione informacje teoretyczne na temat mapowania obiektowo-relacyjnego. W punkcie 3.1 omówione są zasady działania mechanizmu na przykładzie *persistence API* języka Java. Kolejna sekcja (3.2) przedstawia wyniki badania wydajności istniejącej implementacji JPA dla Apache Cassandra. Na podstawie uzyskanych wyników Autor formułuje wnioski zaprezentowane w punkcie 3.3, które są przyczynkiem do przedstawienia autorskiej koncepcji mapowania opisanej w punkcie 3.5.

Rozdział 4 zawiera opis wszystkich unikalnych cech, które zawiera zaproponowane przez Autora narzędzie do modelowania dziedziny w Cassandrze. W sekcji 4.1 wypisane są mechanizmy, dla których wsparcie posiada przygotowany interfejs. Punkt 4.2 wprowadza w podstawowe pojęcia związane z biblioteką. Sekcja 4.3 przedstawia sposób definiowania modelu. Punkt 4.4 to obszernie rozważania na temat modelowania zależności pomiędzy danymi, wraz z kompleksowymi wynikami badań, które pomagają wybrać poprawny model w zależności od scenariusza użycia. W sekcji 4.5 zaprezentowane zostały wszystkie wzorce projektowania schematu, dla których wsparcie udostępnia interfejs modelowania. Punkty 4.6 oraz 4.7 specyfikują zakres obsługi specyficznych mechanizmów silnika bazodanowego. W sekcji 4.8 Autor przedstawia mechanizm migracji danych dedykowany dla Apache Cassandra. Ostatni punkt rozdziału (4.9) przedstawia narzędzia do profilowania aplikacji udostępnione przez system modelowania obiektowego.

Rozdział 2

Apache Cassandra

Apache Cassandra jest bazą danych NoSQL¹, która powstała w wyniku połączenia rozwiązań wykorzystywanych w Dynamo² oraz BigTable³. Cassandra początkowo była rozwijana dla potrzeb portalu społecznościowego Facebook. Baza danych powstała z myślą o rozwiązaniu problemu pełnotekstowego przeszukiwania skrzynek odbiorczych użytkowników, w których dziennie zapisywane były miliardy wiadomości. Głównym celem, do których dążyli twórcy Cassandra była możliwość wykorzystania jej do przechowywania ogromnych ilości danych w bardzo rozproszonym środowisku, gdzie awarie pojedynczych węzłów zdarzają się na porządku dziennym. W tych warunkach baza danych musi zapewniać szybki i niezawodny dostęp do danych. [3]

Apache Cassandra wykorzystywana jest w wielu serwisach na całym świecie. Najbardziej znaczące przykłady użycia produkcyjnego to eBay, Instagram oraz GitHub⁴. Największa światowa instalacja Cassandra obejmuje oko-

¹NoSQL (ang. Not Only SQL) - podzbiór baz danych, które zapewniają inne sposoby modelowania dziedziny niż tradycyjny model oparty na tabelach i relacjach.

²Amazon DynamoDB - zdecentralizowana, wysoce skalowalna baza danych typu klucz-wartość.

³Google BigTable - rozproszony system bazodanowy, który dobrze skaluje się dla ogromnych ilości danych.

⁴eBay, Instagram, Github - przykłady dużych portali internetowych. eBay to serwis aukcyjny, Instagram to portal społecznościowy oparty o publikację zdjęć wykonanych telefonami komórkowymi, a Github to usługa pozwalająca na przechowywanie i wersjonowanie

ło 15000 węzłów, na których przechowywane jest łącznie ponad 4 petabajty danych. [4]

W przeciwieństwie do relacyjnych baz danych, Apache Cassandra nie zapewnia wsparcia dla reguły ACID⁵. Zamiast tego zostały zrealizowane postulaty twierdzenia CAP: „we współdzielonym systemie plików można zachować maksymalnie dwie z trzech właściwości: spójności, dostępności oraz podatności na partycjonowanie”. [5] Apache Cassandra priorytetyzuje właściwości dostępności oraz partycjonowania. Spójność danych jest odwrotnie proporcjonalna i może być regulowana w zależności od czasu odpowiedzi. Wysoka spójność danych oznacza wolniejszą odpowiedź bazy.

2.1 Model danych

Model danych Apache Cassandra jest analogiczny do BigTable. [2] Można przedstawić go jako dwuwymiarowa mapa trójek wartości:

`Map<RowKey, Map<ColumnKey, Triple<Value, Timestamp, TTL>>>`

gdzie `RowKey` to identyfikator wiersza, `ColumnKey` to identyfikator kolumny, `Value` to wartość komórki, `Timestamp` to czas aktualizacji komórki, a `TTL` to czas życia danej wartości. [6] Na rysunku 2.1 przedstawiona jest schematyczna ilustracja wiersza danych. Pogrubiona wartość w lewej komórce to klucz wiersza, natomiast wyróżnione wartości w pierwszym wierszu oznaczają klucze poszczególnych kolumn. Każda komórka składa się z trzech wartości: wartości, czasu życia oraz „odcisku czasu”.

kodu źródłowego aplikacji.

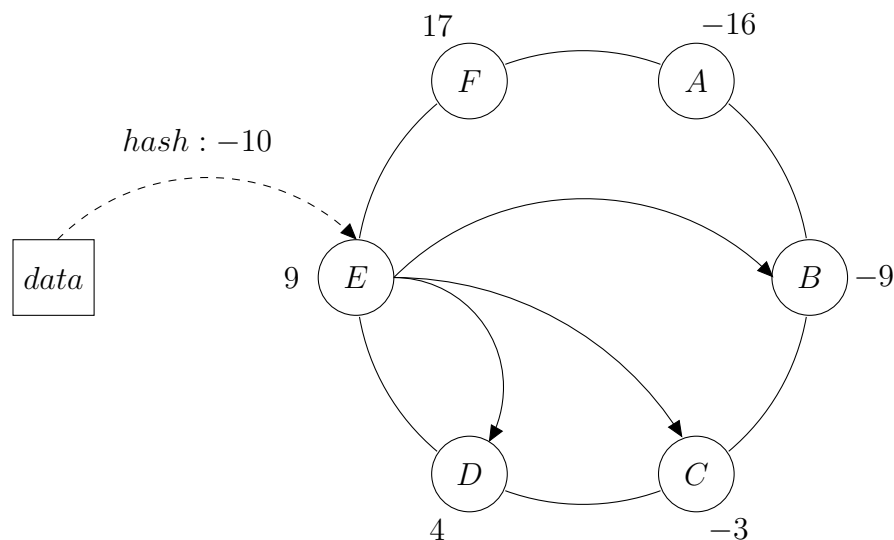
⁵ACID (ang. Atomic, Consistency, Isolation, Durability) - zasada atomowości, spójności, izolacji i trwałości. Wymienione cechy gwarantują poprawne przetwarzanie transakcji w bazach danych.

	ABC	DEF	...	XYZ
123	test value	another test value	...	not a test value
456	20	∞	...	∞
	1291987837942000	1291980736812000	...	1291980736212000

Rysunek 2.1: Przykładowy wiersz modelu danych o identyfikatorze 123456. Wartość komórki (123456, DEF) to „another test value”.

2.2 Dystrybucja danych

Do dystrybucji danych wykorzystywana jest funkcja skrótu, która zachowuje kolejność elementów. Węzły są rozmieszczone w topologii pierścienia. Algorytm dystrybucji zostanie omówiony na przykładzie ze schematu 2.2:



Rysunek 2.2: Schematyczna ilustracja dystrybucji danych w bazie danych Apache Cassandra.

1. Każdemu z węzłów $\{A, B, C, D, E, F\}$ przypisywany jest token, który zawiera się w zakresie wartości przyjmowanych przez funkcję skrótu. Strategię wyboru tokenu można konfigurować. Przykładową strategią

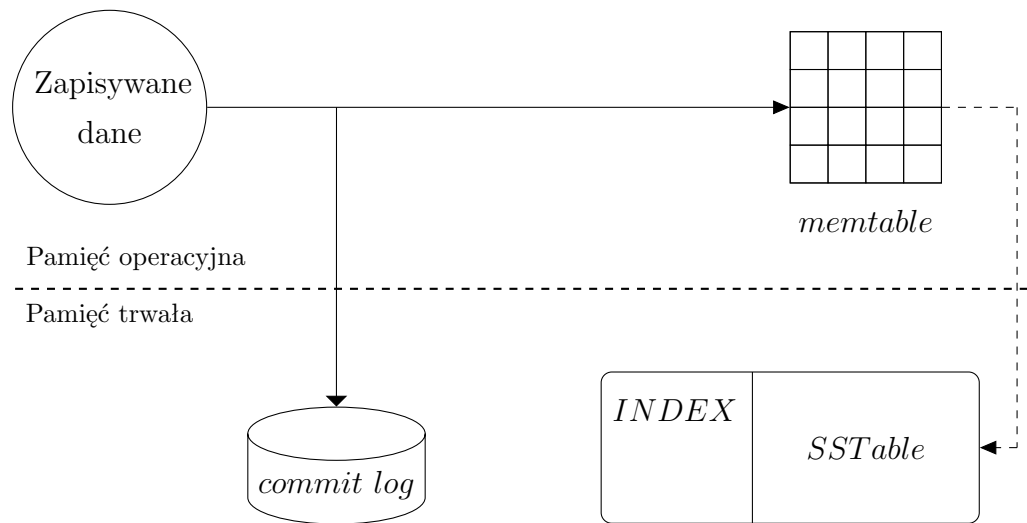
jest wybór losowy. W omawianym przykładzie węzłom zostały przypisane tokeny o wartościach $\{-16, -9, -3, 4, 9, 17\}$.

2. Użytkownik bazy danych przesyła żądanie do dowolnego węzła, który pełni funkcję koordynatora dla danej operacji. Koordynator nadzoruje wpisanie danych do odpowiednich węzłów. W omawianym przykładzie rolę koordynatora pełni węzeł E .
3. Każdy węzeł przechowuje dane, których funkcja skrótu zawiera się w przedziale $(token_{n-1}, token_n]$, gdzie n to numer kolejny węzła rosnący zgodnie z ruchem wskazówek zegara. W przykładzie węzeł C przechowuje wiersze o wartościach funkcji skrótu z przedziału $(-9, -3]$, natomiast węzeł D z przedziału $(-3, 4]$. Wartości funkcji obliczane są cyklicznie, stąd węzeł A przechowuje wiersze o skrócie z przedziału $(-\infty, -16] \cup (17, \infty)$. W przykładzie wiersz o kluczu z funkcją skrótu wartości -10 zostanie utrwalony na węźle B .
4. Dane replikowane są na n węzłach, gdzie n to wartość konfigurowalnego współczynnika replikacji. Poza węzłem macierzystym (wyznaczanym w punkcie 3 algorytmu) dane są replikowane na $n - 1$ kolejnych (zgodnie z ruchem wskazówek zegara) węzłach. W omawianym przykładzie dane zostaną zreplikowane na węzłach C i D .

2.3 Algorytmy zapisu/usuwania danych

Algorytm zapisu danych w Apache Cassandra został schematycznie przedstawiony na diagramie 2.3. [25] Składa się on z następujących kroków:

1. Do węzła przesyłane są dane, które mają zostać na nim zapisane.
2. Dane kopiowane są do dwóch struktur:
 - *memtable* - przechowuje wiersze w pamięci operacyjnej,



Rysunek 2.3: Schematyczna ilustracja algorytmu zapisu danych dla Cassandra.

- *commit log* - przechowuje informacje o kolejnych zapisach wykonywanych do bazy danych.
3. Podstawową strukturą do pobierania danych jest *memtable*. Jest ona umieszczona w szybkiej pamięci, zapewnia więc krótki czas dostępu do danych. Wykorzystanie trwałej struktury *commit log* pozwala na odtworzenie zawartości *memtable* w przypadku awarii wymagającej nagłego restartu maszyny, na przykład braku zasilania. Przy ponownym uruchomieniu węzła Cassandra odtwarza on kolejno wszystkie zapisy, które znajdują się w *commit logu*.
 4. W przypadku przepełnienia *memtable* wykonywane jest „spłukiwanie” zawartości pamięci na dysk do struktury nazywanej *SSTable*. Poza danymi zawiera ona indeks, który pozwala na szybki dostęp do danych wierszy. Jej charakterystyczną cechą jest niezmienność. Po przepisaniu danych z *memtable* do *SSTable* druga z nich nie jest już modyfikowana. Z tego względu wiersze o tym samym kluczu są często podzielone pomiędzy wiele takich struktur. „Spłukiwanie” może zostać uruchomione

manualnie z wykorzystaniem polecenia `nodetool flush`.

5. Ostatnim etapem zapisu danych jest kompresja. Aktualizacja danych nie może nadpisywać istniejących wierszy, bo struktura *SSTable* jest niezmienna. Zamiast tego wstawia ona rekordy z późniejszym „odciskiem czasu”. Powoduje to nadmiarowe zużycie przestrzeni dyskowej. Kompresja pozwala usunąć nieaktualne wiersze poprzez przepisanie istniejących *SSTable* na nowe, uporządkowane struktury. Kompresja może być wywoływana automatycznie w zależności od wybranej strategii lub manualnie, z wykorzystaniem polecenia `nodetool compact`.

Niemodyfikowalność struktur *SSTable* ma swoje konsekwencje także w przypadku usuwania danych. Wiersz nie może zostać fizycznie skasowany. Zamiast tego oznacza się go specjalnym znacznikiem *tombstone*. Wiersze oznaczone tą flagą są usuwane na etapie kompresji danych.

2.4 Obszary zastosowania Cassandra

Specyficzny schemat danych Apache Cassandra i brak wielu mechanizmów znanych z relacyjnych systemów bazodanowych powodują, że nie jest ona najlepszym wyborem do przechowywania uniwersalnych danych:

- Brak wsparcia dla transakcji znacząco utrudnia wykorzystanie Cassandra w dziedzinach, dla których spójność danych jest kwestią kluczową. Przykładem może być obszar finansowy. W systemach relacyjnych spójność jest zapewniona poprzez mechanizmy bazy danych. W przypadku Apache Cassandra odpowiedzialność zostaje przeniesiona na aplikacje dostępowe, co znacząco zwiększa ilość pracy koncepcyjnej i zwiększa ryzyko błędów.
- Brak możliwości złączania powoduje, że model danych musi być projektowany w oparciu o wykonywane do niego odwołania. Utrudnia to

rozbudowę aplikacji. Przykładowo dla relacji jeden do wielu modelowanej poprzez tabelę, w której identyfikatorem wiersza jest identyfikator obiektu nadrzędnego, a w kolumnach wpisywane są identyfikatory obiektów podrzędnych, nie jest możliwe zwrócenie wszystkich obiektów nadrzędnych wskazujących na dany obiekt podrzędny bez modyfikacji schematu.

- Brak złączeń implikuje denormalizację modelu danych. Wprowadza to problemy z zachowaniem spójności. Aby zaktualizować adres, który jest przechowywany w formie zdenormalizowanej w tabeli użytkownika oraz zamówienia należy znaleźć wszystkie rekordy, które wskazują na ten adres. Odpowiedzialność za spójność danych ponownie zostaje przeniesiona na aplikację, która wykorzystuje bazę.

W tabeli 2.4 przedstawiono porównanie różnych aspektów Apache Cassandra i relacyjnych systemów bazodanowych. Analizując informacje zebrane w tabeli i przedstawione wcześniej w rozdziale można wyciągnąć następujące wnioski:

- Modelowanie dziedziny danych jest znacznie prostsze w przypadku systemów relacyjnych. Na etapie projektowania nie trzeba brać pod uwagę wykonywanych odwołań, potencjalnych problemów niespójności, a także możliwych kierunków rozbudowy aplikacji.
- Apache Cassandra najlepiej sprawdza się w zastosowaniu do *big data*. Do poszczególnych fragmentów definicji wielkich zbiorów danych, która została przytoczona w sekcji 1.2, można przyporządkować cechy charakteryzujące Cassandrę:
 - „dużych rozmiarów” - liniowa skalowalność horyzontalna pozwala na obsługę informacji o dowolnej wielkości, pod warunkiem zapewnienia odpowiedniej liczby węzłów,
 - „różnorodnych” - model danych nie jest sztywny, może być dostosowany do zmieniających się danych w trakcie działania aplikacji,

	Apache Cassandra	RDBMS
Skalowalność horyzontalna	Liniowa skalowalność horyzontalna.	Pod warunkiem wykorzystania specjalnych narzędzi i/lub technik, na przykład <i>shardingu</i> ⁶ .
Mechanizmy zachowywania spójności	Brak mechanizmów zachowywania spójności danych.	Mechanizmy i model zapewniające wysoką spójność danych. Wsparcie dla transakcyjności; znormalizowany model.
Odporność na awarie	Wysoka odporność na awarie zapewniona między innymi przez replikację danych pomiędzy węzłami.	Niska odporność na defekty. Pojedynczy punkt awarii.
Optymalizacja operacji	Optymalizacja pod kątem szybkości zapisu.	Optymalizacja pod kątem szybkości wykonywania zapytań.
Koncepcja	Projekt z myślą o opisywaniu szeregu chronologicznego danych. [27]	Projekt do opisu uniwersalnych, ustrukturyzowanych danych.
Przystosowanie modelu danych	Model odpowiedni dla dynamicznych struktur danych zmieniających się w trakcie działania systemu. [26]	Model odpowiedni dla niezmiennych struktur danych.

Rysunek 2.4: Porównanie Apache Cassandra z relacyjnymi bazami danych.

⁶Sharding (od *shard* - ang. kawałek) - technika polegająca na podziale zbioru danych na niezależne partycje w zależności od ich cech, na przykład grupując użytkowników po lokalizacji geograficznej.

- „często zmieniających się” - wbudowana obsługa szeregu chronologicznego oraz optymalizacja czasu zapisu zostały zaprojektowane pod kątem dynamicznie zmieniających się danych.

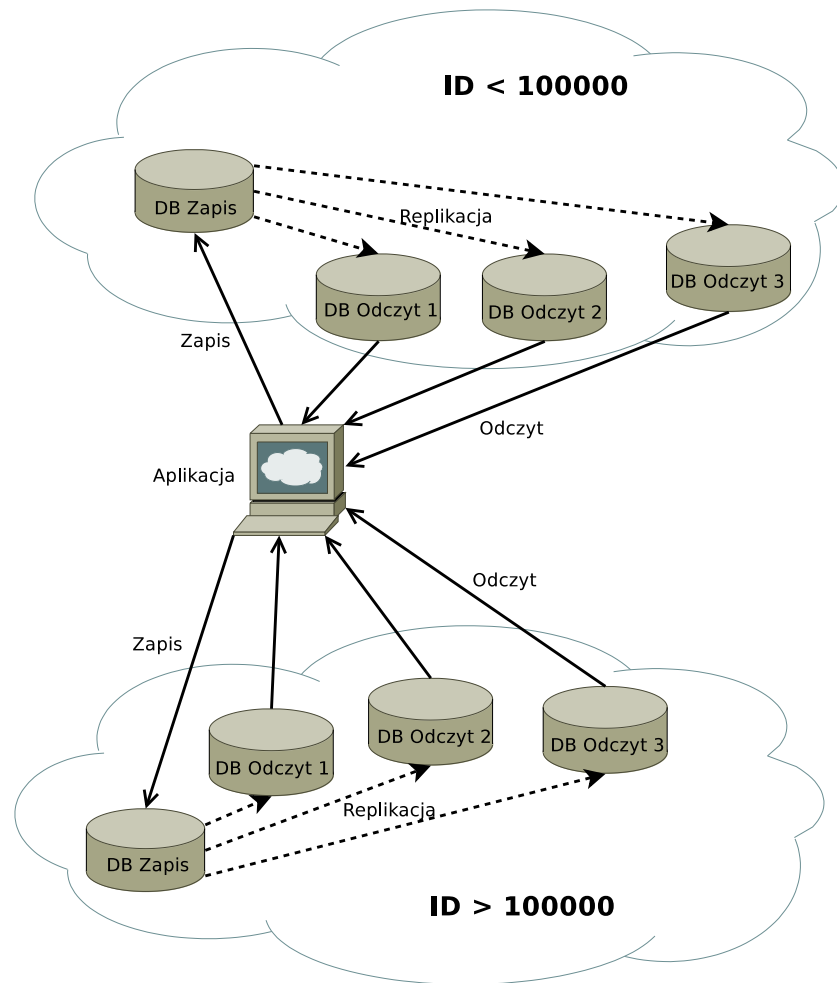
Według twórców Cassandra [29] typowymi obszarami jej zastosowań są:

- Aplikacje wykonujące OLTP⁷, a więc transakcje charakteryzujące się wysoką współbieżnością, krótkim czasem odpowiedzi, małymi zapytaniami oraz pracą na dużych zbiorach danych. [30]
- Zarządzanie danymi uszeregowanymi w czasie.
- Pobieranie i analiza danych z urządzeń o wysokiej częstotliwości odświeżania informacji.
- Aplikacje typu SaaS⁸ oparte o intensywne wykorzystanie usług sieciowych.
- Zarządzanie mediami strumieniowymi (muzyka, filmy) oraz nieustrukturyzowanymi danymi (dobrym przykładem są portale społecznościowe).
- **Systemy „zapisochłonne”.**

Z przedstawionych scenariuszy Autor zdecydował się wyróżnić ostatni. Istotnym zagadnieniem w przypadku wielkich danych jest operacja zapisu, która musi być wykonywana w czasie rzeczywistym. Aby omówić problem należy przedstawić typową architekturę wykorzystywaną dla relacyjnych bazach danych do przechowywania wielkich danych. [32] Przedstawia ją diagram 2.5.

⁷On-Line Transaction Processing - przetwarzanie transakcji sieciowych.

⁸Software as a Service (ang. oprogramowanie jako usługa sieciowa) - aplikacje, które mogą być wykorzystywane przez sieć bez pobierania i instalacji na stacji roboczej. [31] Przykładem takiej aplikacji jest Google Docs (<https://docs.google.com/>), który umożliwia edycję dokumentów, arkuszy kalkulacyjnych oraz prezentacji z wykorzystaniem przeglądarki internetowej.



Rysunek 2.5: Schematyczna ilustracja partycjonowania relacyjnej bazy danych (określanego angielskim terminem *sharding*).

Na diagramie przedstawiono aplikację, która odwołuje się do relacyjnej bazy danych. Została ona podzielona na dwie partycje względem identyfikatora wiersza (mniejszy lub większy od 100000). Każda partycja zawiera jeden węzeł, który wykonuje operacje zapisu (DB Zapis). Do niego odwołuje się aplikacja z żądaniem zapisu danych. W każdej partycji znajdują się również trzy węzły, które pozwalają na odczyt danych (DB Odczyt). Aplikacja zwraca się o odczyt danych kolejno do różnych węzłów. Pozwala to zbalansować obciążenie. Należy zauważyć, że w przedstawionej architekturze aplikacja musi

wiedzieć, do której partycji należy się odwołać dla wyszukiwanych danych.

Proces skalowania systemu o przedstawionej architekturze jest dwustopniowy:

1. Skalowanie operacji odczytu. Może zostać przeprowadzone dla każdej partycji osobno. Polega na dodaniu nowego węzła, który służy do odczytu danych. Jest to stosunkowo proste, należy jednak uwzględnić dodany punkt dostępu w puli serwerów, do których odwołuje się aplikacja.
2. Skalowanie operacji zapisu. W tym przypadku nie wystarczy dodanie drugiego węzła odpowiedzialnego za zapis, gdyż dane i tak muszą zostać zreplikowane pomiędzy węzłami. Z tego względu obciążenie obu serwerów pozostałoby bez zmian. Skalowanie operacji zapisu to skomplikowany proces wymagający stworzenia nowej partycji i modyfikacji reguł przydzielania rekordów do poszczególnych fragmentów.

Podobna architektura wykorzystywana jest w komercyjnych rozwiązaniach, takich jak **MySQL Cluster CGE**⁹. Występują tam dodatkowe warstwy pośredniczące, które ukrywają szczegóły wykonywania zapytań przed aplikacją. Sama zasada działania pozostaje jednak analogiczna.

Problemem w tej architekturze pozostaje skalowanie operacji zapisu. Może ono zostać wykonane poprawnie tylko pod warunkiem, że istnieje funkcja partycjonująca dla danego natężenia operacji zapisu, która dzieli całość danych na dostatecznie wydajne (a więc małe) fragmenty. Okazuje się, że w praktyce nie zawsze jest to możliwe. Właśnie ten problem przyczynił się do powstania projektu Apache Cassandra w Facebooku.

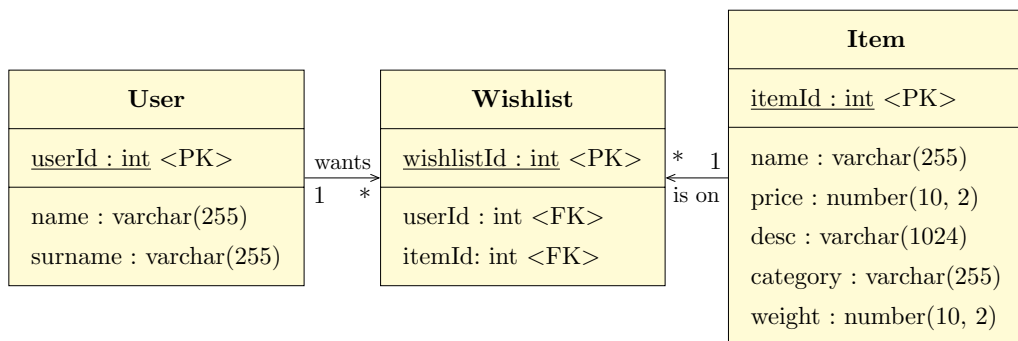
Wszystkie przypadki użycia wymienione przez twórców Cassandra pokrywają scenariusze, w których ich system bazodanowy sprawdza się lepiej niż podejście relacyjne w subiektywnym odczuciu. Omówiony szerzej przypadek z zapisem jest jednak jedynym, który z przyczyn technologicznych może być nierealizowalny w świecie SQL-owych baz danych.

⁹Strona domowa projektu - <http://www.mysql.com/products/cluster/>.

2.5 Struktura danych a modelowanie

Struktura i mechanizm dystrybucji danych wykorzystywany w Apache Cassandra zmienia podejście do modelowania dziedziny znane z relacyjnych baz danych. Zbudowanie efektywnego modelu danych Cassandra wymaga skupienia się w podobnym stopniu na zdefiniowaniu encji z modelowanego świata, jak również na analizie odwołań, które będą wykonywane do obiektów z tego świata. [7]

Założmy, że celem jest modelowanie danych dla sklepu internetowego. Zakupów dokonują użytkownicy, którzy mogą wstawić wiele przedmiotów z oferty sklepu na listę życzeń. W przypadku baz opartych o język SQL jest to klasyczny problem relacji typu wiele-do-wielu, do modelowania których wykorzystywana jest najczęściej tabela pośrednia.



Rysunek 2.6: Modelowanie listy życzeń w relacyjnej bazie danych.

Diagram prezentujący zamodelowaną relację dla listy życzeń jest przedstawiony na rysunku 2.6. W tabeli Użytkownik (*User*) przechowywane są imię, nazwisko oraz identyfikator. W tabeli Przedmiot (*Item*) znajduje się nazwa, cena, a także inne właściwości: opis, kategoria oraz waga. Tabela Lista życzeń (*Wishlist*) łączy ze sobą użytkownika i przedmiot poprzez wykorzystanie kluczy obcych.

Powyższy model jest wykorzystywany w widoku listy życzeń na profilu użytkownika. Na liście życzeń prezentowane są informacje o nazwie przedmiotu oraz jego cenie. Po kliknięciu nazwy użytkownik przenoszony jest do

strony przedmiotu. Na listingu 2.7 zaprezentowano zapytanie, które wyświetla listę życzeń.

```
SELECT item.name, item.price
FROM Item item, Wishlist wishlist
WHERE wishlist.userId = 202;
```

Rysunek 2.7: Zapytanie, które pobiera wszystkie przedmioty na liście życzeń użytkownika o identyfikatorze 202.

Cassandra umożliwia utworzenie dokładnej repliki relacyjnego modelu danych. Zostało to przedstawione na rysunku 2.8.

User		name	surname	Wishlist		userId	itemId
	123	Janusz	Kowalski		51	123	579
		name	surname			userId	itemId
	124	Marcin	Nowak		52	124	232

Item		name	price	desc	category	weight
	232	Master Chef	20.34	Recipes	BOOKS	0.2
		name	price	desc	category	weight
	579	Seat Hit	159.99	Armchair	FURNITURE	10.8

Rysunek 2.8: Wynik błędnego przeniesienia relacyjnego modelu danych do Cassandra.

Taki model jest jednak niepoprawny. Nie umożliwia on filtrowania zawartości listy życzeń po identyfikatorze użytkownika. Wynika to z faktu, że pobranie odpowiednich wierszy listy życzeń wymaga znajomości ich identyfikatorów, podczas gdy widok dysponuje wyłącznie odniesieniem do użytkownika. Błąd ten można łatwo naprawić zastępując encję *Wishlist* encją *WishlistByUser*, co przedstawia diagram 2.9.

Poprawiony model można poddać dalszej optymalizacji. Wyświetlenie listy życzeń użytkownika wymaga odwołania do encji *Item*, w której znajdują

WishlistByUser		579
	123	-
		232
	124	-

Rysunek 2.9: Definicja encji listy życzeń umożliwiająca filtrowanie względem użytkownika.

się informacje o nazwie i cenie przedmiotu. Ponieważ przedmioty mogą być rozmieszczone na różnych węzłach, silnik Cassandra nie może wykonać złączenia w sposób optymalny - zapytanie o każdą pozycję listy życzeń jest wykonywane osobno. W celu przyspieszenia wykonywania operacji należy wykonać denormalizację. Dołączając do encji *WishlistByUser* informacje o nazwie i cenie produktu można uniknąć wykonywania kosztownych złączeń. Pozostałe dane przedmiotu zostaną pobrane dopiero po przejściu na jego stronę. Efekt denormalizacji jest przedstawiony na diagramie 2.10.

WishlistByUser		579
	123	(„Master Chef”, 20.43)
		232
	124	(„Seat Hit”, 159.99)

Rysunek 2.10: Zdenormalizowana postać listy życzeń.

2.6 CQL

Efektywne modelowanie i obsługa danych w Apache Cassandra wymaga dobrej znajomości wewnętrznej struktury bazy danych. Dodatkowym utrudnieniem przy korzystaniu z początkowych wersji Cassandra była konieczność wykorzystania skomplikowanego interfejsu programistycznego opartego

o wywoływanie zdalnych procedur Thrift¹⁰. Thrift jest platformą pozwalającą budować aplikacje przenośne między różnymi językami programowania. Dzięki temu rozwiązaniu baza danych dostępna była dla wszystkich platform. Niestety, skutkiem ubocznym było skomplikowanie interfejsu dostępowego.

Wraz z wydaniem 1.2 Apache Cassandra wprowadzony został nowy interfejs dostępu do tej bazy danych. Interfejs ten nosi nazwę CQL¹¹ i jest językiem zapytań, którego składnia wzorowana jest na SQL. Poza podobieństwami składniowymi języki te nie mają cech wspólnych. Nie są wzajemnie zgodne. W chwili obecnej CQL jest rekomendowanym standardem komunikacji z Apache Cassandra. [8] Na listingu 2.11 przedstawiono zapytanie w języku CQL, które opisuje omawianą wcześniej encję *User*. Wynikiem wykonania tego zapytania jest prosty schemat modelu danych zaprezentowany na diagramie 2.12.

```
CREATE TABLE User (
    userId uuid PRIMARY KEY,
    name text,
    surname text);
```

Rysunek 2.11: Zapytanie CQL, które tworzy encję *User*.

User		name	surname
	userId	null	null

Rysunek 2.12: Wynik wykonania zapytania 2.11.

2.7 Przykład modelowania dziedziny danych

W Internecie znajdują się przykładowe modele danych stworzone przez użytkowników Apache Cassandra. Zapewniają one materiał do nauki dla nowych

¹⁰Dokumentacja interfejsu dostępna jest pod adresem <https://wiki.apache.org/cassandra/API10>.

¹¹CQL (ang. Cassandra Query Language) - język zapytań Cassandra.

użytkowników bazy. Prezentują sprawdzone sposoby rozwiązywania zadanego problemu koncentrując się na omówieniu decyzji projektowych, które należy podjąć aby zbudować efektywny schemat. Proces tworzenia wydajnych modeli danych zostanie przez Autora przedstawiony z wykorzystaniem takiego przykładu. Do modelowania zostanie wykorzystany język CQL, gdyż począwszy od wersji CQL3 jest on oficjalnym standardem zalecanym przez autorów Cassandra.

2.7.1 Twissandra

Projekt Twissandra został stworzony przez Tylera Hobbsa i opublikowany pod adresem <https://github.com/twissandra/twissandra>. Według informacji na stronie „Twissandra jest przykładowym projektem stworzonym do nauki i zaprezentowania jak używać Cassandra. Po uruchomieniu projektu zaprezentowana zostanie strona internetowa, która posiada funkcjonalność podobną do serwisu Twitter”. [33]

Serwis Twitter¹² pozwala na publikację krótkich wiadomości tekstowych o długości maksymalnie 160 znaków. Jego funkcjonalność jest wzorowana na SMS¹³. Różnica polega na tym, że wiadomość nie jest kierowana do konkretnego adresata (adresatów), ale publicznie dostępna dla wszystkich odwiedzających profil danego użytkownika.

Głównym widokiem w Twitterze jest oś czasu. Znajdują się na niej wpisy ułożone w odwrotnym porządku chronologicznym. Istnieją trzy typy osi czasu:

Publiczna przedstawia wpisy wszystkich użytkowników.

Użytkownika przedstawia wpisy danego użytkownika oraz wszystkich śledzonych przez niego osób.

Profilowa przedstawia wyłącznie wpisy danego użytkownika.

¹²Dostępny pod adresem <https://twitter.com>.

¹³Short Message Service (ang. usługa krótkich wiadomości) - usługa pozwalająca na przesyłanie pomiędzy telefonami komórkowymi krótkich (do 160 znaków) wiadomości. [34]

Pierwszym krokiem jest zdefiniowanie modelu użytkownika. Dla uproszczenia encja posiadać będzie wyłącznie dwa parametry - nazwę oraz hasło - oba zapisywane jako tekst. Istotną kwestią jest wybór klucza głównego. W przypadku relacyjnych baz danych identyfikatory najczęściej tworzone są w sposób sztuczny, a odpowiednie wiersze wybierane są z zastosowaniem zapytań. W przypadku Cassandra wybór klucza głównego jest dużo bardziej istotny:

- Cassandra nie posiada mechanizmu generacji unikatowego klucza głównego. Zamiast tego wykorzystywane są typy `uuid` lub `timeuuid`, [35] które są obliczane przez aplikację kliencką.
- Ze względu na architekturę Cassandra nie potrafi efektywnie pobrać wszystkich wierszy należących do danej tabeli. Oznacza to, że bez używania indeksów (co wiąże się z tworzeniem kolejnych tabel) **jedyną możliwością odwołania się do wiersza jest podanie jego identyfikatora**. Wybór losowo/sekwencyjnie generowanego klucza może sprawić, że pobieranie danego wiersza będzie bardzo utrudnione.

Dla modelu użytkownika dobrym identyfikatorem jest nazwa (zakładając jej unikatowość). Jest to dana, która będzie znana w dowolnym kontekście odwołania. Na listingu 2.13 przedstawiono zapytanie, które pozwala utworzyć i wstawić użytkownika do bazy danych. W wyniku tego zapytania zostanie utworzona tabela o strukturze przedstawionej na diagramie 2.14.

```
CREATE TABLE users (  
    username text PRIMARY KEY,  
    password text)  
INSERT INTO users (username, password) VALUES (  
    'jturek',  
    'UnsafePassword')
```

Rysunek 2.13: Tworzenie i wstawianie użytkownika w Twissandrze.

jturek	password
	UnsafePassword

Rysunek 2.14: Poglądowy wynik struktury uzyskanej po wykonaniu zapytań 2.13.

Kolejnym krokiem jest umożliwienie śledzenia użytkowników. W relacyjnym modelu danych do przechowywania takiej informacji można użyć tabeli **followers**, która posiada trzy kolumny: sztuczny klucz główny, identyfikator śledzonego użytkownika oraz identyfikator śledzącego. Ta sama architektura w Cassandra będzie błędna, gdyż nie da się zadać zapytania, które zwróci wszystkich użytkowników śledzących profil o zadanej nazwie:

- Nie ma możliwości efektywnego pobrania wszystkich wpisów z tabeli **followers**, a więc nie da się ich filtrować po nazwie śledzonego użytkownika.
- Istnieje możliwość stworzenia indeksu na kolumnie z identyfikatorem śledzonego użytkownika. Dzięki temu możliwe będzie wybranie zadanych wpisów w tabeli **followers**.
- Indeks tworzy niewidoczną tabelę, w której kluczem głównym wiersza jest jedna z unikalnych wartości komórki (w podanym przypadku identyfikator śledzonego użytkownika), a w kolumnach przechowywane są identyfikatory wszystkich wierszy, które posiadają taką wartość komórki.

Analizując opis metody tworzenia indeksów łatwo dojść do wniosku, że można przekształcić model tabeli **followers** w taki sposób, aby ich używanie nie było potrzebne. Przyjmując identyfikator śledzonego użytkownika jako klucz wiersza, a informację o nazwach profili śledzących przechowując w kolejnych kolumnach, możliwe jest pobranie kompletu informacji w jednym zapytaniu.

CQL3 umożliwia wykorzystanie dwóch typów kluczy do dwuwymiarowego podziału wierszy [36]:

Partition key określa, na którym węźle znajduje się dany wiersz. W niskopoziomowym schemacie danych jest implementowany na kluczu wiersza.

Clustering key określa kolumny, względem których porządkowane są wiersze dla danej partycji. W niskopoziomowym schemacie danych jest implementowany na kolumnach.

Dzięki wykorzystaniu dwóch rodzajów kluczy można zdefiniować efektywny schemat dla tabeli **followers**. Odpowiednie zapytanie zostało przedstawione na listingu 2.15. Poglądowy wygląd fizycznego schematu danych uzyskanego dla omawianego zapytania został zaprezentowany na diagramie 2.16.

```
CREATE TABLE followers (  
    username text,  
    follower text,  
    since timestamp,  
    PRIMARY KEY (username, follower))  
INSERT INTO followers (username, follower, since) VALUES (  
    'jturek',  
    'jkowalski',  
    1412020800)  
INSERT INTO followers (username, follower, since) VALUES (  
    'jturek',  
    'mnowak',  
    1409601600)
```

Rysunek 2.15: Śledzenie użytkowników w Twissandrze.

W zaprezentowanym przykładzie kolumna **username** pełni rolę klucza **partition**, natomiast **follower** to klucz typu **clustering**. Dla powyższego przypadku można wypisać wszystkie wiersze tabeli **followers** dla użytkownika o kluczu **jturek**. Odpowiednie zapytanie zostało zaprezentowane na listingu 2.17. Wynik tego zapytania przedstawia tabela 2.18.

jturek	follower:jkowalski	follower:mnowak
	1412020800	1409601600

Rysunek 2.16: Poglądowy wynik struktury uzyskanej po wykonaniu zapytań 2.15.

```
SELECT * FROM followers WHERE username = 'jturek'
```

Rysunek 2.17: Zapytanie, które wybiera wszystkie elementy tabeli *follower* dla użytkownika *jturek*.

username	follower	since
jturek	jkowalski	1412020800
jturek	mnowak	1409601600

Rysunek 2.18: Wybór wszystkich wpisów tabeli *follower* dla użytkownika *jturek*.

Usunięcie kolumny **follower** z klucza utworzonej tabeli spowodowałoby, że dany użytkownik mógłby być śledzony tylko przez jedną osobę. Z kolei przesunięcie tej kolumny do klucza typu **partition** pozwoliłoby pobierać wiersze wyłącznie przy podaniu kombinacji wartości **username** oraz **follower**, co nie jest intencją omawianego modelu.

Tabela **followers** umożliwia reprezentację relacji śledzenia tylko w jedną stronę. Do realizacji widoku, na którym będą wyświetlone wszystkie wpisy osób śledzonych przez użytkownika wymagana jest informacja odwrotna. Autor Twissandry sugeruje stworzenie komplementarnej tabeli **friends**, w której zapisane będą nazwy osób śledzonych przez dany profil. W omawianym przypadku równie wydajnym rozwiązaniem jest wykorzystanie indeksów zaimplementowanych w CQL. Stworzenie indeksu na kolumnie **follower** w tabeli **followers** przy pomocy zapytania przedstawionego na listingu 2.19 umożliwia wyszukanie wszystkich użytkowników śledzonych przez danego. Przykładowe zapytanie realizujące takie wyszukanie zostało zaprezentowane na listingu 2.20.

```
CREATE INDEX ON followers (follower)
```

Rysunek 2.19: Zapytanie, które tworzy indeks na kolumnie *follower* w tabeli *followers*.

```
SELECT * FROM followers WHERE follower = 'jkowalski'
```

Rysunek 2.20: Zapytanie, które wybiera wszystkie osoby śledzone przez użytkownika *jkowalski* dzięki indeksowi.

Kolejnym elementem modelu jest definicja tabeli, która umożliwia przechowywanie wpisów. Pełni ona rolę normalizacyjną i nie będzie wykorzystywana w widokach osi czasu. W tym przypadku można posłużyć się sztucznym identyfikatorem, gdyż pobieranie wierszy następować będzie wyłącznie w kontekście znanego klucza wpisu. Zapytania pozwalające stworzyć strukturę tabeli oraz wstawić do niej przykładowy wpis zostały zaprezentowane na listingu 2.21. Poglądową strukturę fizyczną uzyskaną w wyniku ich wykonania prezentuje natomiast diagram 2.22.

```
CREATE TABLE tweets (
    tweet_id uuid PRIMARY KEY,
    user text,
    body text)
INSERT INTO tweets (tweet_id, user, body) VALUES (
    f8122d00-2f99-11e4-babf-0002a5d5c51b,
    'jturek',
    'Sample')
```

Rysunek 2.21: Tabela wpisów w Twissandrze.

	user	body
f8122d00-2f99-11e4-babf-0002a5d5c51b	'jturek'	'Sample'

Rysunek 2.22: Poglądowy wynik struktury uzyskanej po wykonaniu zapytań 2.21.

Tabela `tweets` nie może być wykorzystana do konstrukcji widoków osi czasu. Aby wyświetlić wszystkie wpisy użytkownika i śledzonych przez niego osób należałoby złączyć trzy tabele - `tweets`, `users` oraz `followers`. Apache Cassandra nie umożliwia jednak wykonywania złączeń. Możliwe jest manualne scalanie tabel w aplikacji, ale podobnie jak dla relacyjnych baz danych jest to rozwiązanie skrajnie niewygodne. W celu uniknięcia złączeń należy zbudować widok osi czasu w osobnej tabeli. Jej zawartość będzie uzupełniana podczas dodawania nowych wpisów przez użytkowników. Na osi czasu wyświetlane będą dwie informacje - nazwa użytkownika publikującego oraz treść wpisu.

Autor Twissandry sugeruje, aby w tabeli osi czasu umieścić trzy informacje - nazwę użytkownika, a także czas publikacji oraz identyfikator wpisu. Przy tak zdefiniowanym schemacie na etapie budowy widoku dla każdego wpisu osobno pobierana jest jego treść z tabeli `tweets`. Nie jest to kluczowy problem wydajnościowy, jednakże dzięki denormalizacji dodatkowe pobranie może być łatwo wyeliminowane. Wystarczy uwzględnić treść wpisu w tabeli osi czasu. Zapytanie, które tworzy schemat `timeline` zostało zaprezentowane na listingu 2.23. Poglądowy kształt wynikowej struktury danych przedstawia diagram 2.24.

W zapytaniu tworzącym tabelę `timeline` wykorzystano klauzulę sortowania `WITH CLUSTERING ORDER BY`. Pozwala ona na układanie kolejności wierszy względem wartości klucza typu `clustering`. W przykładzie 2.23 sortowanie wykonywane jest względem wartości kolumny `time`, w porządku malejącym¹⁴. Oznacza to, że wiersze pobrane dla danej nazwy użytkownika będą posortowane w odwrotnym porządku chronologicznym (od najnowszych).

Należy zauważyć, że w Cassandrze możliwe jest sortowanie wierszy tylko i wyłącznie po wartości klucza typu `clustering`. Wynika to z reprezentacji wewnętrznej danych. Klucz typu `partition` jest tłumaczony na identyfikator wiersza w Cassandrze. Wiersze są rozdzielane w architekturze pierścienia, w zależności od wartości skrótu identyfikatora. Zmiana ich sortowania nie

¹⁴Świadczy o tym słowo kluczowe *DESC* - od angielskiego *descending* - malejąco.

```

CREATE TABLE timeline (
    username text,
    time timeuuid,
    tweet_id uuid,
    body text,
    PRIMARY KEY (username, time)
) WITH CLUSTERING ORDER BY (time DESC)
INSERT INTO timeline (username, time, tweet_id, body) VALUES (
    'jturek',
    915d6810-308c-11e4-95de-f9043d8d556b,
    f8122d00-2f99-11e4-babf-0002a5d5c51b,
    'Sample')

```

Rysunek 2.23: Oś czasu w Twissandrze z modyfikacją wprowadzającą denormalizację treści wpisu.

'jturek'	time: 915d6810-308c-11e4-95de-f9043d8d556b	
	tweet_id	body
	f8122d00-2f99-11e4-babf-0002a5d5c51b	'Sample'

Rysunek 2.24: Poglądowy wynik struktury uzyskanej po wykonaniu zapytań 2.23.

jest możliwa, gdyż zaprzecza to pryncypiom architektonicznym Cassandra. Sortowanie klucza typu `clustering` jest możliwe i polega na szeregowaniu kolumn w wewnętrznej reprezentacji danych.

W powyższym schemacie kluczowym problemem jest reprezentacja publicznej osi czasu. Autor Twissandry wykorzystuje do tego celu zdefiniowaną wcześniej tabelę `timeline` używając specjalnego identyfikatora użytkownika (`!PUBLIC!`). Rozwiązanie to nie jest jednak poprawne:

- Cassandra może przechowywać maksymalnie miliard kolumn w jednym wierszu (przy założeniu, że cały węzeł wypełnia pojedynczy wiersz). [37] Oznacza to, że w tabeli `timeline` dla klucza `!PUBLIC!` można umieścić

ograniczoną liczbę wpisów. Istnieje realne ryzyko, że liczba ta okaże się zbyt mała.

- Architektura w istotny sposób degradowe dobrą skalowalność horyzontalną, którą cechuje się Cassandra. Wiersz z publiczną osią czasu będzie znacznie dłuższy niż przeciętny wiersz z osią czasu użytkownika, gdyż znajdują się w nim wszystkie wpisy utworzone w serwisie. W zaprezentowanym schemacie publiczna oś czasu będzie umieszczona w jednym fizycznym wierszu Cassandra¹⁵, a więc węzły przechowujące ją będą nieproporcjonalnie obciążone w porównaniu do pozostałych.

Rozwiązaniem tego problemu jest wprowadzenie większej granulacji publicznej osi czasu. Można to osiągnąć rozbudowując klucz typu `partition` o kolejną składową. Sugerowanym przez Autora rozbiciem osi czasu jest podział względem dnia publikacji wpisu. Przykładowo aplikacja może odwoływać się do klucza `!PUBLIC!2014-07-01` dla wpisów opublikowanych pierwszego lipca 2014 roku, natomiast do pobrania danych z kolejnego dnia wykorzystywać identyfikator `!PUBLIC!2014-07-02`.

Dużym ułatwieniem jest fakt, że aplikacja nie musi samodzielnie budować i zarządzać skomplikowanymi identyfikatorami. Język CQL umożliwia wykorzystanie złożonego klucza typu `partition`. [38] Listing 2.25 prezentuje w jaki sposób zdefiniować oś czasu z dodatkowym podziałem względem dnia publikacji. Zmiana obejmuje dodanie kolumny tekstowej `day` oraz uwzględnienie jej w klauzuli `PRIMARY KEY`. W fizycznym modelu danych wpisy o jednakowych wartościach nazwy użytkownika i dnia zostaną pogrupowane w wiersze, natomiast wpisy dla tej samej osoby, ale pochodzące z różnych dni będą rozdzielone pomiędzy węzłami.

Wykorzystując klucze złożone należy pamiętać, że istotna jest kolejność nawiasów zastosowanych w klauzuli `PRIMARY KEY`. W przypadku pominięcia wewnętrznego nawiasu dla pary (`username`, `day`) nowa kolumna zostałaby częścią klucza typu `clustering`, który byłby opisany dwójką (`day`, `time`).

¹⁵Należy pamiętać, że wiersze w fizycznym modelu danych są niepodzielne.

```
CREATE TABLE timeline (  
    username text,  
    day text,  
    time timeuuid,  
    tweet_id uuid,  
    body text,  
    PRIMARY KEY ((username, day), time)  
) WITH CLUSTERING ORDER BY (time DESC)  
INSERT INTO timeline (username, day, time,  
    tweet_id, body) VALUES (  
    'jturek',  
    '2014-08-01',  
    915d6810-308c-11e4-95de-f9043d8d556b,  
    f8122d00-2f99-11e4-babf-0002a5d5c51b,  
    'Sample')
```

Rysunek 2.25: Oś czasu w Twissandrze partycjonowana względem dnia utworzenia wpisu.

2.8 Modelowanie - wzorce i antywzorce

Pomimo, że CQL znacząco upraszcza modelowanie w Cassandrze stworzenie efektywnego schematu danych nie jest zadaniem prostym. Aby ułatwić ten proces twórcy i użytkownicy Cassandra zaczęli opisywać wzorce i antywzorce modelowania oraz dostępu do danych. Pełnią one analogiczną rolę do wzorców i antywzorców projektowych znanych z inżynierii oprogramowania. Na ich opis składa się możliwie ogólna definicja problemu, a także poprawny (lub niepoprawny) sposób jego rozwiązania.

Przykładem antywzorca modelowania jest kolejka oparta na kolumnach. [9] Kolejka to struktura danych, w której ilości wykonywanych operacji wstawiania, usuwania i odczytu są do siebie zbliżone. W przypadku Cassandra usuwanie kolumn z wiersza nie jest wykonywane natychmiast po odebraniu

żądania. Zamiast tego usunięta kolumna jest oznaczana specjalnym znacznikiem *tombstone* i fizycznie usuwana dopiero po upływie pewnego czasu. Takie działanie przyspiesza znacząco operację usuwania kosztem operacji odczytu. Kiedy w wierszu występuje wiele znaczników *tombstone* operacje filtrowania zakresu kolumn wykonywane są znacznie wolniej.

Przykładem wzorca dostępu do danych jest unikanie konfliktów synchronizacji poprzez uaktualnianie wyłącznie zmodyfikowanych wartości. [10] Entcja *Item* z diagramu 2.8 ma 5 właściwości. Wyświetlenie ekranu aktualizacji przedmiotu wymaga pobrania zawartości całego wiersza z bazy danych. Wzorec stanowi, że jeżeli na tym ekranie zostanie zmieniony wyłącznie opis to do Cassandra należy przesłać żądanie uaktualniające zawartość wyłącznie jednej komórki - *desc*. Pozostałe wartości z formularza mogą być już nieaktualne. Przesłanie kompletu informacji poskutkowałoby nadpisaniem aktualnych wartości. Postępowanie według wzorca pozwala wykorzystywać mechanizm rozwiązywania konfliktów wbudowany w Cassandrę.

Rozdział 3

Mapowanie obiektowo-relacyjne

Mapowanie obiektowo-relacyjne (zwyczajowo określane skrótem ORM od angielskiego terminu *object-relational mapping*) to technologia, która pozwala automatyzować połączenie pomiędzy relacyjnym modelem baz danych, a paradygmatem programowania zorientowanego obiektowo. [18] W szerokiej perspektywie mapowanie obiektowo-relacyjne może być postrzegane jako próba przełożenia tabelarycznej reprezentacji danych umieszczonej w pamięci masowej na dowolną reprezentację danych w pamięci operacyjnej. [19]

Potrzeba stworzenia mechanizmów ORM to naturalna konsekwencja powstawania aplikacji zorientowanych na przetwarzanie dużej ilości współdzielonych danych. Obsługa dostępu do bazy danych, czyli nieodłączna część współczesnych aplikacji biznesowych, wymaga napisania dużych ilości kodu źródłowego. Kod ten jest powtarzalny pomiędzy poszczególnymi rodzajami danych i aplikacjami. Jego stworzenie jest pracochłonne, a nie stanowi żadnej wartości funkcjonalnej dla aplikacji. Mechanizmy ORM w sposób znaczący redukują ilość kodu niezbędnego do komunikacji z systemami bazodanowymi. Ich główne zalety to:

- Możliwość definiowania modelu danych poprzez tworzenie klas reprezentujących obiekty wykorzystywane w aplikacji.
- Implementacja typowych operacji na danych: tworzenia, aktualizacji, usuwania i wyszukiwania.

- Eliminacja konieczności lub uproszczenie zarządzania połączeniami, sesjami i transakcjami bazodanowymi.
- Zwiększenie bezpieczeństwa aplikacji. Mechanizmy ORM dostarczają narzędzi do ochrony przed atakami typu SQL Injection¹. [11]
- Uniwersalne wsparcie dla wielu silników bazodanowych, które mogą różnić się od siebie implementowanymi standardami języka SQL.
- Umożliwienie wymiany środowiska bazodanowego bez konieczności modyfikacji kodu aplikacji.

O powszechnym użyciu mechanizmów ORM w projektach informatycznych świadczą statystyki:

- W Internecie dostępne są tysiące implementacji mechanizmów ORM o otwartych źródłach dla dziesiątek różnych języków programowania. Wyszukiwanie frazy *object-relational mapping* w witrynie <http://github.com> zwraca około 45 tysięcy wyników.
- Popularnym narzędziem ORM jest biblioteka Hibernate dla języka Java. Wyszukiwanie frazy `<artifactId>hibernate-core</artifactId>`² w witrynie <http://github.com> zwraca około 151 tysięcy wyników.

O popularności mechanizmów ORM ponadto fakt powstawania oficjalnych standardów mapowania obiektowo-relacyjnego włączanych do specyfikacji języków programowania. Przykładem takiego standardu jest JPA³ dla języka Java.

Mapowanie obiektowo-relacyjne nie jest pozbawione wad. Opanowanie podstaw mechanizmów ORM jest często trudniejsze niż nauka języka SQL.

¹SQL Injection (ang. wstrzykiwanie SQL) - typ ataku polegający na wykorzystywaniu specjalnie spreparowanych wartości, które dołączone do szablonu zapytania SQL powodują szkodliwe działania niezgodne z intencją programisty.

²Jest to fraza będąca częścią deklaracji biblioteki Hibernate w popularnym narzędziu do budowania projektów w języku Java - Maven.

³Java Persistence API - szerszy opis znajduje się w sekcji 3.1.

Wynika to ze złożoności takich mechanizmów. Przykładowo na projekt Hibernate ORM przypada ponad milion linii kodu źródłowego. Inną często przytaczaną wadą mapowania obiektowo-relacyjnego jest drastyczny spadek wydajności wynikający z natury mechanizmów ORM. Prowadzi to do sytuacji, w których organizacje decydują się na stworzenie własnego, dedykowanego dla danego problemu oprogramowania. Ilość projektów, które wykorzystują systemy mapowania obiektowo-relacyjnego pozwala jednak twierdzić, że wady te są akceptowalne w obliczu licznych zalet ORM.

3.1 Java Persistence API

Java Persistence API to interfejs programistyczny, który ma za zadanie uprościć tworzenie, zarządzanie i zapisywanie obiektów, które reprezentują dane z baz relacyjnych. [20] Interfejs JPA operuje na POJO⁴, które dekorowane są adnotacjami umożliwiającymi konfigurację reguł mapowania. Przykład takiego obiektu prezentuje listing 3.1. Adnotacja `@Entity` określa klasę, która reprezentuje encję danych. Parametr `@Id` dekoruje identyfikator encji, natomiast `@GeneratedValue` oznacza, że będzie on ustawiany automatycznie. Adnotacje `@Column` informują, że dane zmienne będą mapowane na kolumny bazodanowe. Parametr `@ManyToMany` definiuje relację wiele-do-wielu.

Definicja encji 3.1 nie specyfikuje kompletnego przykładu. Brakuje przede wszystkim wartości konfiguracyjnych, takich jak nazwa tabeli, z której encja będzie mapowana, nazwy kolumn czy też sposób reprezentacji relacji wiele-do-wielu. W rzeczywistości JPA nie jest tak czytelne jak sugeruje przedstawiony przykład. Parametry konfiguracyjne zajmują większą część definicji klasy i zdarza się, że najważniejsza informacja - specyfikacja samego obiektu, który reprezentuje dane - jest mocno przysłonięta.

Na listingu 3.2 przedstawiono przykładowy kod, który służy do zapisy-

⁴Plain Old Java Object (ang. dosłownie „prosty, stary obiekt Java”) - termin określający zwykły obiekt języka Java, używany jako przeciwieństwo Enterprise Java Bean, czyli specjalnego obiektu, który stosował się do wielu ściśle określonych reguł.

```
@Entity
public class User {
    @Id
    @GeneratedValue
    private Integer id;
    @Column
    private String name;
    @Column
    private String surname;
    @Column
    private String password;
    @ManyToMany
    private List<Item> wishlistItems;
    // getters & setters
}
```

Rysunek 3.1: Przykładowy obiekt specyfikujący użytkownika w standardzie JPA.

wania użytkownika w bazie danych. Do obsługi operacji zapisu wykorzystywany jest zarządca encji (`EntityManager`). Za jego pomocą rozpoczynana jest transakcja (`beginTransaction()`), wywoływane jest żądanie utrwalenia obiektu (`persist()`), a następnie transakcja jest wykonywana (`commit()`).

```
User user = new User();
// setting field values
EntityManager manager = managerFactory.createEntityManager();
manager.getTransaction().beginTransaction();
manager.persist(user);
manager.getTransaction().commit();
manager.close();
```

Rysunek 3.2: Przykładowy kod zapisujący użytkownika w bazie danych w standardzie JPA.

Do najbardziej popularnych implementacji Java Persistence API należą Hibernate ORM, OpenJPA oraz Toplink. Dużą przeszkodą jest fakt, że interfejs JPA nie jest wyczerpująco określony. Opisuje on tylko najbardziej pod-

stawowe operacje na obiektach bazodanowych. Z jednej strony pozostawia to twórcom bibliotek dużą elastyczność w kwestii implementacji. Z drugiej sprawia, że w przypadku wykonywania operacji bardziej skomplikowanych niż zapisanie i odczytanie mapowanego obiektu, programista często jest zmuszony sięgać do elementów specyficznych dla danej implementacji JPA. Sprawia to, że większość aplikacji wykorzystujących wymienione biblioteki jest nieprzenośna pomiędzy mechanizmami.

3.2 Kundera

Wraz z pojawieniem się języka CQL pojawiła się szansa wykorzystania istniejących implementacji mechanizmów ORM do przechowywania danych z użyciem silnika Cassandra. Motywacją do stworzenia takiego rozwiązania jest zwiększenie wydajności istniejących aplikacji bez modyfikacji ich kodu. Minimalna rekonfiguracja aplikacji i podłączenie jej pod klaster bazodanowy skutkowałyby zwiększeniem wydajności działania aplikacji. Ponadto wykorzystanie istniejących interfejsów mogłoby umożliwić obsługę różnych typów baz danych należących do ruchu NoSQL.

Przykładem projektu, który wykorzystuje mechanizm ORM do obsługi NoSQLowych baz danych jest Kundera.[12] Kundera zapewnia implementację mapowania obiektowego zgodną ze standardem JPA 2.0 między innymi dla Cassandra, HBase, MongoDB oraz Neo4j. Dodatkowo biblioteka wspiera obsługę wielu mechanizmów równocześnie.

3.2.1 Wydajność

Wyniki pomiarów na stronie domowej projektu Kundera świadczą o tym, że biblioteka nie wprowadza znacznych narzutów wydajnościowych względem bezpośredniego wykorzystania interfejsu bazy danych. Ważniejsze jest jednak sprawdzenie jak duży narzut wprowadza dostosowanie relacyjnego modelu danych do Apache Cassandra. W tym celu Autor przeprowadził testy porównawcze masowego wstawiania i pobierania obiektów. Dla każdej

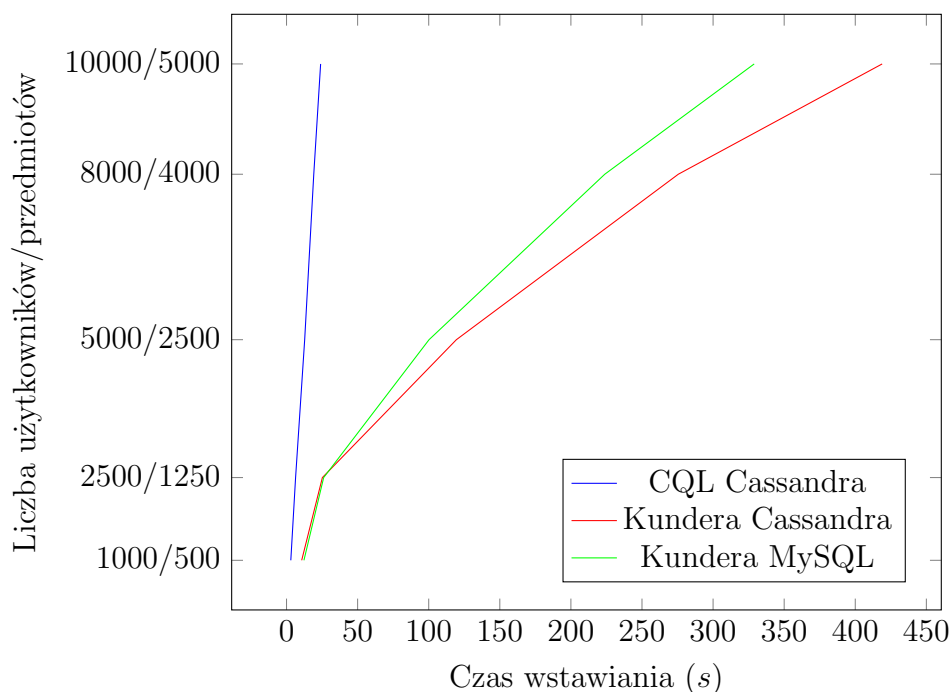
z dwóch operacji zostały przeprowadzone trzy testy. Pierwszy test mierzy czas referencyjny. Jest wykonywany dla zdenormalizowanego modelu Cassandra przedstawionego na diagramie 2.10. Drugi test sprawdza czasy dla znormalizowanego modelu zaprezentowanego na diagramie 2.6, opisanego w JPA i uruchomiony na Cassandrze z wykorzystaniem biblioteki Kundera. Trzeci test wykorzystuje identyczny opis modelu jak drugi, jednakże wykonywany jest dla silnika MySQL. Wszystkie testy przeprowadzane były jednowątkowo. Czas opóźnień przesyłania danych jest pomijalny, gdyż testy były prowadzone na środowisku lokalnym.

W trakcie testów mierzony był całkowity czas wykonania danej operacji. Przyjęte zostały następujące założenia:

- Liczba użytkowników i przedmiotów są parametryzowane i skalowane liniowo.
- Liczba przedmiotów, które użytkownik może mieć na swojej liście życzeń zawiera się w przedziale $[0; 10]$.
- Odczytywana jest parametryzowalna liczba użytkowników.

Wyniki czasu wstawiania wielu rekordów zostały przedstawione na wykresie 3.3. Czasy pobierania wielu rekordów zostały zebrane na wykresie 3.4.

Wyniki zebrane na wykresie 3.3 pokazują, że wykorzystanie mechanizmów mapowania obiektowo-relacyjnego dla bazy danych Cassandra jest bardzo złym wyborem. Narzut związany z konwersją modelu danych do postaci relacyjnej jest tak duży, że w praktyce pojedynczy węzeł Cassandra osiąga gorsze wyniki zapisu niż baza danych MySQL. Zmiana silnika bazodanowego dla istniejącego kodu nie tylko nie poprawi wyników wydajnościowych, ale może wręcz spowodować spowolnienie działania aplikacji. Jedynym zyskiem z takiego rozwiązania będzie możliwość wykorzystania natywnego mechanizmu klastrowania węzłów Cassandra. Dzięki temu wyeliminowany zostanie pojedynczy punkt awarii systemu. Ten sam efekt można jednak uzyskać stosując



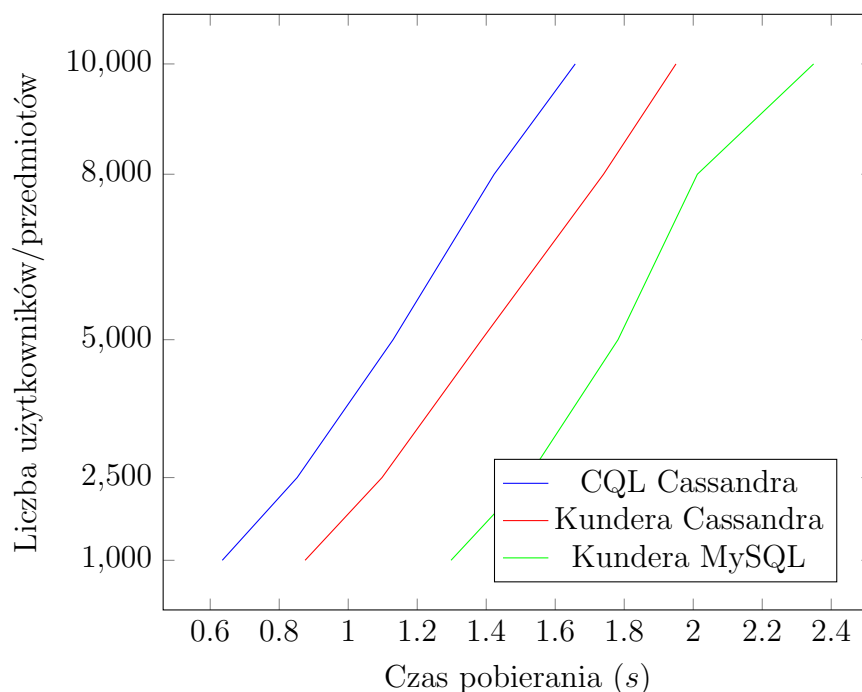
Rysunek 3.3: Porównanie czasu wstawiania wielu rekordów.

rozwiązana dedykowane dla relacyjnych baz danych. Przykładami takich rozwiązań mogą być MySQL Cluster CGE dla bazy danych MySQL oraz Oracle RAC dla baz Oracle Database Enterprise Edition.

Potencjalnym zastosowaniem mapowania obiektowo-relacyjnego dla bazy danych Cassandra są środowiska integracyjne dla wielu aplikacji. Wykorzystując homogeniczny model danych można odwoływać się do różnych silników bazodanowych. W bibliotece Kundera istnieje takie rozwiązanie. Zostało nazwane Polyglot Persistence⁵. W praktyce dostosowanie modelu danych ORM do istniejących schematów baz danych jest problematyczne.

Wykres 3.3 demonstruje jak ogromną przewagę szybkości przy zapisie posiada poprawnie zaprojektowany model danych Apache Cassandra. Przypuszczalnie osiągnięty czas mógłby być jeszcze lepszy, gdyż ograniczenie wy-

⁵Polyglot Persistence (dosłowne tłumaczenie to „zapis poliglotyczny”) - opis mechanizmu znajduje się na stronie <https://github.com/impetus-opensource/Kundera/wiki/Polyglot-Persistence>.



Rysunek 3.4: Porównanie czasu pobierania wielu rekordów.

dajności zapisu następowało prawdopodobnie po stronie klienta testowego. Zgodnie ze specyfikacją Apache Cassandra jest w stanie obsługiwać bez opóźnień znacznie więcej jednoczesnych żądań zapisu.

Czasy pobierania listy życzeń użytkownika przedstawione na wykresie 3.4 nie różnią się znacząco od siebie. Przewaga szybkości pobierania danych z wykorzystaniem CQL ponownie wynika z zastosowania lepszego modelu danych. Denormalizacja pól encji *Item* pozwala pominąć pobieranie dla każdego użytkownika dodatkowego wiersza z bazy danych.

3.3 Wnioski z badania wydajności

Na podstawie wyników przeprowadzonych badań można wnioskować, że mechanizmy mapowania obiektowo-relacyjnego nie nadają się do wykorzystania przez bazę danych Apache Cassandra. Wynika to z faktu, że mechanizmy te wspomagają normalizację i zachowywanie relacji między encjami, natomiast

efektywne modelowanie dziedziny danych w Cassandra dąży do denormalizacji i maksymalnego uproszczenia zależności pomiędzy obiektami.

Drastyczna różnica wydajności pomiędzy szybkością zapisu danych w modelach zoptymalizowanym i niezoptymalizowanym świadczy o tym, że do efektywnego modelowania dziedziny w Apache Cassandra niezbędna jest wiedza na temat sposobu przechowywania danych przez tę bazę i podejmowanie świadomych wyborów. Dostarczona przez ORM-y mechanizmy mogą natomiast wyłącznie spowolnić korzystanie z Cassandra.

Podczas badania wydajności Autor natrafił również na problemy z przenośnością kodu źródłowego pomiędzy różnymi systemami bazodanowymi. Wykonanie operacji wstawiania partii danych w bazie relacyjnej oraz Cassandra wymagało modyfikacji kodu źródłowego aplikacji testowej. W relacyjnej bazie danych nie można było ukończyć testu wstawiania i pobierania danych wykorzystując jedynie elementy zdefiniowane w JPA. Pomimo wielokrotnej zmiany konfiguracji w pewnym momencie testu dochodziło do całkowitego zablokowania puli połączeń. Dopiero przejście do transakcji bezstanowej zdefiniowanej w Hibernate ORM⁶ umożliwiło zakończenie testu.

Opisane problemy z przenośnością kodu mogą być następstwem fundamentalnych różnic pomiędzy relacyjnymi systemami bazodanowymi a Apache Cassandra. Przykładowo transakcja, czyli jedno z podstawowych pojęć związanych z zapytaniami w języku SQL, nie istnieje w systemie Cassandra. Dostosowanie implementacji do interfejsu mapowania obiektowo-relacyjnego wymaga sztucznego blokowania zapytań po stronie klienta, co prowadzi do wystąpienia niepotrzebnych opóźnień i nadmiernego wykorzystania mechanizmów synchronizacji wielowątkowej. Z drugiej strony brak lub nieodpowiednie wykorzystanie transakcji w przypadku bazy relacyjnej prowadzi do błędów zapytań.

Wyniki badań nie świadczą o niemożliwości implementacji mapowania

⁶Hibernate ORM - środowisko aplikacyjne implementujące standard JPA dla baz relacyjnych. Zdefiniowana w ramach tego środowiska transakcja bezstanowa nie jest częścią interfejsu JPA, a więc nie może być wspierana przez Kunderę.

obiekтового dla Cassandra. Dla optymalnej wydajności należy jednak zastosować rozwiązanie dedykowane:

- Definiowanie modelu powinno być skoncentrowane wokół wewnętrznej reprezentacji danych w Cassandrze. W stosunku do mapowania obiektowo-relacyjnego kluczowe są aspekty obniżenia istotności relacji i wsparcie dla operacji denormalizacji z poziomu interfejsu.
- Mapowanie powinno w prosty i transparentny sposób udostępniać realizację poprawnych wzorców modelowania. Przykładowo lista wartości powinna móc być zamodelowana przynajmniej na dwa sposoby: jako kolumna typu list lub jako lista kolumn o nazwach zawierających wartości listy. Mapowanie obiektowe dla Cassandra powinno umożliwiać wymienienie fizycznej struktury danych bez zmiany wykorzystującego ją kodu źródłowego.
- Mapowanie obiektowe dla Cassandra powinno być proste. Język CQL, w porównaniu do SQL, ma prostą składnię i nie udostępnia wielu operacji. Ponadto skomplikowany mechanizm mapowania mógłby niekorzystnie kontrastować z wydajną bazą danych. Przykładem nieczytelnego mapowania z dużą liczbą opcji konfiguracyjnych może być pole encji opisane w JPA przedstawione na listingu 3.5.

3.4 Hibernate OGM

Alternatywnym rozwiązaniem, które dostarcza implementacji Java Persistence API dla baz NoSQL jest Hibernate Object/Grid Mapper (OGM). Hibernate OGM wykorzystuje silnik projektu ORM o tej samej nazwie. Niestety, aktualne wydanie (wersja **4.1.0.Beta5**) wspiera wyłącznie silniki Infispan, Ehcache, MongoDB oraz Neo4j.

Rozwiązanie to jest jednak interesujące z punktu widzenia dalszego rozwoju pracy. Według mapy przyszłych wydań Hibernate OGM w wersji 4.2 wprowadzi wsparcie dla silnika bazy danych Apache Cassandra. [21]

```

@ManyToMany
@JoinTable(name = "wishlist",
    joinColumns = {
        @JoinColumn(name = "userId",
            referencedColumnName = "userId") },
    inverseJoinColumns = {
        { @JoinColumn(name = "itemId",
            referencedColumnName = "itemId") },
        @JoinColumn(name = "price",
            referencedColumnName = "price") }
    },
    foreignKey = @ForeignKey(name = "userId_fk"),
    inverseForeignKey = @ForeignKey(name = "itemId_fk"))
private Set<Item> wishlistItems = new HashSet<Item>();

```

Rysunek 3.5: Przykład nieczytelnego mapowania obiektowo-relacyjnego.

3.5 Koncepcja mapowania dla Cassandra

W większości rozwiązań ORM encja modelu danych jest opisywana jako definicja klasy, której obiekty reprezentować będą instancję tej encji. Rozważmy listing 3.6, który przedstawia hipotetyczne dostosowanie standardu JPA do denormalizacji dla modelu 2.10.

```

@Table(name = "wishlist")
class Wishlist {
    @Id(type = IdType.PARTITION_KEY)
    private String userId;
    @Denormalize(clustering_keys = { "itemId" },
        fields = { "name", "price" })
    private Item item;
}

```

Rysunek 3.6: Hipotetyczny przykład dostosowania definicji JPA do denormalizacji.

W powyższym przykładzie denormalizacja jest modelowana jako zdegenerowany przypadek relacji jeden-do-jednego. Przypisanie zmiennej *item* do obiektu klasy *Wishlist* i utrwalenie tego obiektu w bazie danych spowoduje uzupełnienie wartości kolumn *itemId*, *name* oraz *price* w tabeli *wishlist*.

Po pobraniu obiektu klasy *Wishlist* z bazy danych przechowywany reprezentant klasy *Item* posiada tylko część informacji. Takie rozwiązanie posiada istotną wadę. Podczas korzystania z mechanizmu nie można jednoznacznie stwierdzić czy brak uzupełnienia pola *desc* w obiekcie klasy *Item* jest spowodowany tym, że opisu faktycznie brakuje, czy też operacje dokonywane są na niekompletnym (zdenormalizowanym) obiekcie. Rozwiązaniem pozbawionym tej wady jest denormalizacja jawna przedstawiona na listingu 3.7.

```
@Table(name = "wishlist")
class Wishlist {
    @Id(type = IdType.PARTITION_KEY)
    private String userId;
    @Id(type = IdType.CLUSTERING_KEY)
    private String itemId;
    @Column(name = "name")
    private String name;
    @Column(name = "price")
    private String price;
}
```

Rysunek 3.7: Hipotetyczny przykład zastosowania denormalizacji jawnej w JPA.

Denormalizacja jawna posiada jednak inne wady. Przede wszystkim model nie podaje, że kolumny *itemId*, *name* oraz *price* powinny pochodzić od instancji klasy *Item*. Ponadto użytkownik mechanizmu musi pamiętać o ręcznym wypełnieniu pól, co jest z kolei narażone na błędy. Eleganckie połączenie tych dwóch opcji nie jest możliwe w języku Java.

Ze względu na mechanizm metaklas w Pythonie stworzono eleganckie i proste implementacje ORM dla relacyjnych baz danych. Dwa najpopularniejsze mechanizmy to Django ORM⁷ oraz SQLAlchemy⁸. Wykorzystując

⁷Django ORM - część platformy Django, która odpowiada za operacje bazodanowe. Dokumentacja jest dostępna pod adresem <https://docs.djangoproject.com/en/dev/topics/db/>.

⁸Strona domowa projektu jest dostępna pod adresem <http://www.sqlalchemy.org>.

metaklasy osiągalne jest stworzenie mapowania o interfejsie wymienionych mechanizmów, który rozwiązuje problem zasygnalizowany w przypadku języka Java. Ilustruje to listing 3.8.

```
class Wishlist(Model):
    userId = TextField(partition_key=True)
    item = DenormalizedField(related=Item,
                            clustering_keys=['itemId'],
                            fields=['name', 'price'])

w = Wishlist()
w.item = Item(id='1', name='n', price='10.0')
id = w.item_itemId    # id equals '1'
name = w.item_name     # name equals 'n'
price = w.item_price   # price equals '10.0'
i = w.item              # exception raised
```

Rysunek 3.8: Hipotetyczny przykład denormalizacji w mapowaniu w języku Python.

Dzięki metaklasom można dokonać następujących modyfikacji modelu *Wishlist*:

- Dla każdej instancji *DenormalizedField* automatycznie dodać do klasy pola, które będą przechowywać zdenormalizowane wartości. W przykładzie definicja klasy została rozszerzona o pola *item_itemId*, *item_name*, *item_price*.
- Każdą instancję *DenormalizedField* zamienić na właściwość tylko do zapisu, która automatycznie ustawia wartości zdenormalizowanych pól.

Dzięki takim modyfikacjom możliwe jest ustawianie wartości zdenormalizowanych pól poprzez podawanie całej instancji obiektu *Item*, jak w listingu 3.6. Jednocześnie pobieranie wartości z obiektu możliwe jest tylko poprzez odwołania do zdenormalizowanych pól. Nie da się pomylić ze sobą zdenormalizowanych i pełnych instancji klasy *Item*.

Rozdział 4

Modelowanie obiektowe dla Cassandry

Na podstawie wyników wydajnościowych przedstawionych w sekcji 3.3 Autor zaproponował odmienną koncepcję mechanizmu mapowania obiektowego dla bazy danych Cassandra. Aby uwypuklić różnice pomiędzy przedstawioną propozycją i tradycyjnymi mechanizmami ORM Autor postuluje stosowanie odmiennego nazewnictwa. Omówiona w sekcji 3.5 koncepcja, która będzie rozwijana w dalszej części pracy, nazywana będzie **modelowaniem obiektowym**:

- Mechanizmy mapowania obiektowo-relacyjnego nie wymuszają na użytkowniku kolejności projektowania komponentów. Osiągalne są zarówno podejście *code first*¹, jak i *database first*². W proponowanym przez Autora modelowaniu obiektowym, ze względu na występowanie struktur wysokiego poziomu odpowiadających wzorcom projektowym, podejście *database first* jest możliwe tylko teoretycznie lub w bardzo ograniczo-

¹Code first (ang. dosłownie „najpierw kod źródłowy”) - w kontekście ORM oznacza to modelowanie dziedziny za pomocą kodu źródłowego, z którego następnie generowany jest schemat bazy danych.

²Database first (ang. dosłownie „najpierw baza danych”) - w kontekście ORM oznacza to modelowanie dziedziny jako schemat bazodanowy, do którego następnie pisany jest mapujący kod źródłowy.

nym zakresie. Znacznie bardziej efektywne jest podejście *code first*. Posiada ono wsparcie implementacyjne - generację schematu - i jest rekomendowane przez Autora.

- W mechanizmach mapowania obiektowo-relacyjnego mapowanie odpowiada najczęściej typom danych, natomiast w obrębie danego typu jest ono jednoznaczne. Przykładowo identyfikator konwertowany jest na liczbowy klucz główny, lista przechowywana jest jako relacja, a ciąg znaków przekształcany jest na typ *VARCHAR*. W przypadku modelowania obiektowego dla Cassandra wybór wyznaczać będzie natomiast sposób przechowywania danej wartości. W zależności od narzuconego sposobu modelowania listy będzie mogła zostać ona zrzutowana na wartość kolumny o typie *list*, nazwy kolumn w wierszu lub osobną tabelę.
- Mapowanie obiektowo-relacyjne dostarcza przede wszystkim niskopoziomowe odpowiedniki typów bazodanowych. Wyjątkiem jest relacja wiele-do-wielu. Modelowanie obiektowe dla Cassandra skupia się zarówno na polach niskiego poziomu (przykładowo kolumna o wartości tekstowej), jak i na reużywalności komponentów wysokopoziomowych odpowiadających wzorcom projektowym.
- Częścią modelowania obiektowego dla Cassandra są opcjonalne mechanizmy mapowania obiektowo-relacyjnego. Należą do nich generacja schematu oraz migracje pomiędzy wersjami modelu.

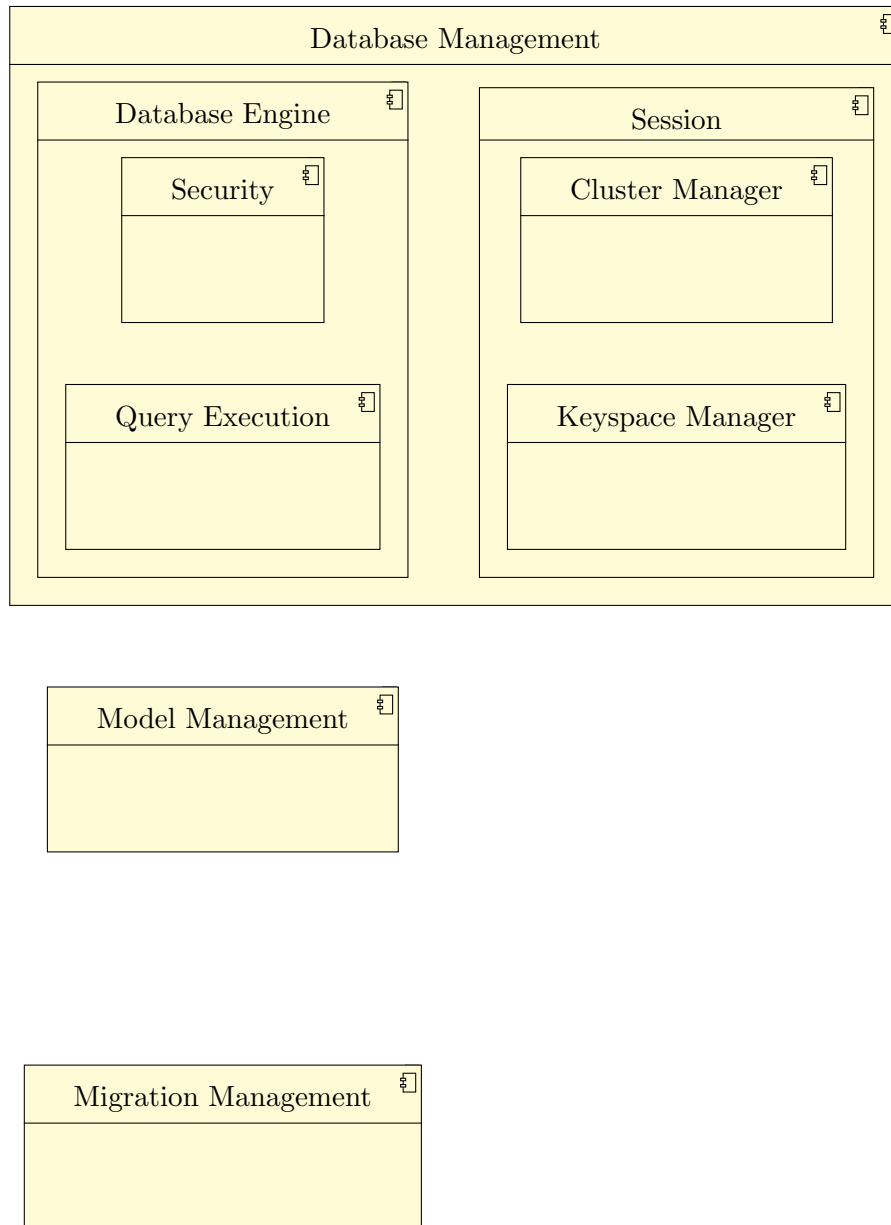
W dalszej części pracy omówiona zostanie konkretna implementacja modelowania obiektowego dostarczona wykonana przez Autora. Pełna nazwa tego projektu to **Object Modeling for Cassandra**. W dalszej części pracy używana będzie nazwa skrócona - *OMC*.

4.1 Object Modeling for Cassandra

OMC jest biblioteką dla języka Python, która dostarcza narzędzia do modelowania dziedziny danych dedykowane dla bazy Apache Cassandra. Pakiet udostępnia moduły, które umożliwiają zarządzanie środowiskiem bazodanowym i pracę z danymi bez odwoływania się do interfejsu Thrift/CQL. Funkcjonalności pakietu obejmują:

- Mapowanie obiektowe dla modelu danych Cassandra:
 - Pobieranie, wyszukiwanie i usuwanie danych za pomocą metod interfejsu.
 - Automatyczne pakowanie i rozpakowywanie wartości do/z postaci obiektów.
- Wsparcie dla specyficznych mechanizmów Cassandra:
 - Indeksy drugiego poziomu.
 - Liczniki.
 - Operacje na partiach danych.
- Wspomaganie modelowania zależności między obiektami z wykorzystaniem denormalizacji.
- Wspomaganie tworzenia modelu danych poprzez dostarczenie/wsparcie reużywalnych struktur wysokiego poziomu:
 - szeregów zdarzeń (ang. *time series*),
 - indeksów wartości unikalnych,
 - kolejek.
- Tworzenie i walidacja zgodności modelu danych z wymaganiami mapowania.
- Automatyczna migracja danych pomiędzy różnymi wersjami modelu.

- Narzędzia do populacji danych testowych i profilowania aplikacji.



Rysunek 4.1: Diagram komponentów mechanizmu Object Modeling for Cassandra.

4.2 OMC - podstawowe pojęcia

Środowisko aplikacyjne OMC definiuje trzy podstawowe pojęcia: model, pole oraz silnik.

4.2.1 Model

Jest to klasa, która opisuje atomowy względem kodu źródłowego obiekt modelujący bazę danych. Obiekt ten, w zależności od definicji, może być mapowany na jedną lub wiele tabel. Jedynym wymaganiem, które musi spełnić klasa należąca do modelu jest dziedziczenie po klasie bazowej `Model`. Ciało modelu stanowią pola.

4.2.2 Pole

Pole jest to zmienna klasy definiującej model. Opisuje atomową wartość z punktu widzenia kodu źródłowego. W zależności od typu pole może być mapowane na jedną kolumnę, wiele kolumn lub autonomiczną tabelę w bazie danych. Aby pole klasy mogło zostać zakwalifikowane jako pole modelu, w definicji klasy należy mu przypisać obiekt dziedziczący po klasie bazowej `Field`.

4.2.3 Silnik

Silnik jest to obiekt, który odpowiada za nawiązywanie połączenia z bazą danych. Po utworzeniu silnika wystarczy przypisać go do całego modelu. Wszystkie operacje bazodanowe będą domyślnie wykonywane przez przypisany silnik.

Utworzenie silnika wymaga podania adresu URL opisującego połączenie, nazwanego z języka angielskiego *connection string*. Strukturę tego adresu przedstawia listing 4.2.

Wszystkie elementy ujęte w znaki `{ }` w adresie 4.2 należy zamienić na wartości według poniższego objaśnienia:

```
cassandra://{ip}:{port}/{keyspace}?rf={rf}&strategy={strategy}
```

Rysunek 4.2: Struktura adresu URL opisującego połączenie z bazą danych Cassandra.

ip adres IP węzła koordynatora, z którym będzie komunikować się kod źródłowy,

port port IP koordynatora (podanie portu jest opcjonalne; wartość domyślna to 9042),

keyspace przestrzeń nazw Cassandra, w której operuje OCM,

rf współczynnik replikacji³ podawany przy tworzeniu przestrzeni nazw, jeżeli ta nie istnieje (parametr jest opcjonalny),

strategy nazwa strategii replikacji podawanej przy tworzeniu przestrzeni nazw, jeśli ta nie istnieje (parametr jest opcjonalny).

Do utworzenia silnika należy wykorzystać metodę klasy **Engine** o nagłówku `create_engine(connection_string)`. Aby przypisać silnik do modelu należy użyć metody `Model.bind(engine)`. Przykład wywołania przypisującego modelowi silnik dla serwera lokalnego i przestrzeni nazw `test` został zaprezentowany na listingu 4.3.

```
engine = Engine.create_engine('cassandra://127.0.0.1/test')
Model.bind(engine)
```

Rysunek 4.3: Przykład wywołania tworzącego silnik i przypisującego go do modelu.

4.3 Definiowanie modelu

Definiowanie modelu odbywa się poprzez stworzenie nowej klasy, która dziedziczy po klasie bazowej **Model**. Podstawowy model opisuje zawartość jednej

³Współczynnik replikacji (ang. replication factor) - określa na ile węzłów kopiowany jest każdy wstawiony wiersz.

tabeli, jego pola zaś odpowiadają kolumnom tej tabeli. Przykładowy model został przedstawiony na listingu 4.4.

```
class SampleTableModel(Model):  
    id = UuidField(type=Uuid, partitioning_key=True)  
    sample_field = TextField()
```

Rysunek 4.4: Przykładowy model danych OMC.

Po zdefiniowaniu klasy modelu oraz utworzeniu i przypisaniu silnika jak na listingu 4.3 uruchomi się algorytm weryfikacji schematu bazodanowego. Algorytm działa następująco:

1. Sprawdzenie czy jest zdefiniowana przestrzeń nazw (dla przykładu: `test`).
2. Jeżeli przestrzeń nazw nie jest zdefiniowana zostanie ona utworzona.
3. Dla każdego zdefiniowanego modelu sprawdź czy istnieje odpowiadająca mu tabela.
4. Jeżeli tabela istnieje, sprawdź czy definicja tabeli odpowiada modelowi. Jeśli nie, zgłoś wyjątek modelu danych.
5. Jeżeli tabela nie istnieje zostanie utworzona.

W wyniku działania algorytmu zostanie utworzona tabela z wykorzystaniem zapytania przedstawionego na listingu 4.5. Należy zauważyć, że OCM w schemacie bazodanowym wykorzystuje notację z podkreśleniami⁴. Gdyby klucz główny nie został zdefiniowany jawnie, mechanizm utworzyłby go sztucznie jako autogenerowane pole o nazwie `id` i typie `timeuuid`.

```
CREATE TABLE sample_table_model (id PRIMARY KEY, sample_field text);
```

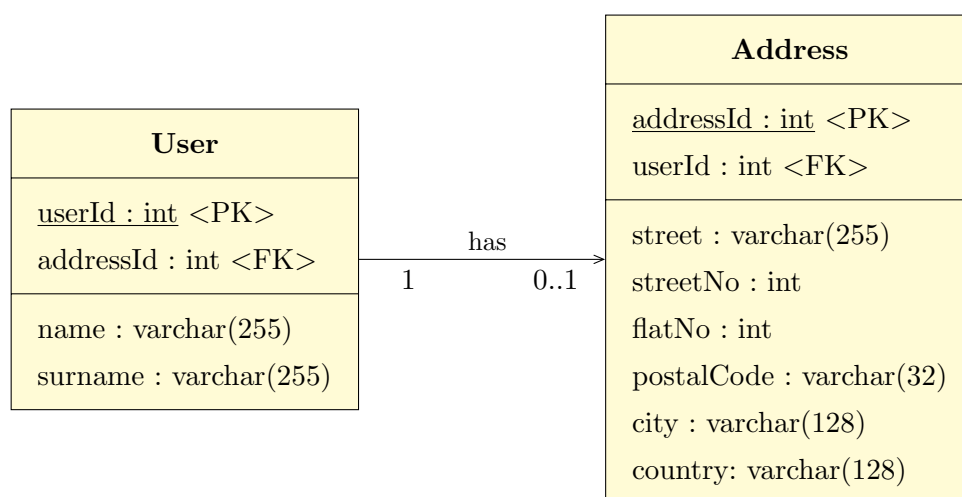
Rysunek 4.5: Zapytanie tworzące tabelę dla modelu 4.4.

⁴Underscore notation - zapis, w którym używane są wyłącznie małe litery, a poszczególne wyrazy są oddzielone znakiem `_`, na przykład: `underscore_notation_example`.

4.4 Modelowanie zależności między danymi

Autor celowo nie używa terminu „relacja”, aby uniknąć skojarzenia z RDBMS⁵. Jak wykazały badania przeprowadzone w sekcji 3.2.1 tradycyjne podejście relacyjne w Apache Cassandra jest bardzo nieefektywne i powinno być unikane.

W przypadku relacyjnych baz danych modelowanie wymaga określenia definicji modelowanych obiektów - encji oraz zależności między nimi. Na diagramie 4.6 zaprezentowano efekt modelowania prostej dziedziny danych, która umożliwia przechowywanie adresu dla użytkownika. Charakterystyka RDBMS pozwala na pobranie danych użytkownika i adresu w jednym odwołaniu do bazy danych.



Rysunek 4.6: Modelowanie adresu użytkownika w relacyjnej bazie danych.

Model ten można dokładnie odwzorować w Apache Cassandra, jednakże jest to rozwiązanie niewydatne. Charakterystyka bazy sprawia, że sprawdzenie kraju zamieszkania użytkownika wymaga wykonania dwóch osobnych zapytań: pobrania użytkownika, a następnie wskazywanego przez odpowiadający mu wiersz adresu. Czas przeznaczony na wykonanie jednego zapytania można aproksymować wzorem:

⁵Relative Database Management Systems - relacyjne systemy zarządzania bazami danych

$$t_{SELECT} = t_{AC} + t_{HASH} + t_{CN} + t_{GET} + t'_{CN} + t'_{AC} \quad (4.1)$$

gdzie t_{AC} , t'_{AC} to czas komunikacji (i komunikacji zwrotnej) pomiędzy aplikacją a koordynatorem, t_{CN} , t'_{CN} to czas komunikacji (i komunikacji zwrotnej) pomiędzy koordynatorem a węzłem przechowującym wyszukiwany wiersz, t_{HASH} to czas obliczania funkcji skrótu dla klucza danego wiersza, a t_{GET} to właściwy czas pobrania wiersza. Pomiedzy czynnikami zachodzi zależność:

$$t_{AC}, t'_{AC}, t_{CN}, t'_{CN} \gg t_{HASH}, t_{GET}. \quad (4.2)$$

Wynika stąd, że wykonanie dodatkowego zapytania wprowadza opóźnienia związane głównie z połączeniem pomiędzy poszczególnymi maszynami. Opóźnienia te w większości przypadków są relatywnie niskie (rzędu kilku milisekund), więc nie jest to krytyczny problem wydajnościowy. Ponieważ można ich uniknąć, należy traktować je jako błąd projektowania modelu.

Ze względu na specyfikę modelu danych Apache Cassandra poprawnym podejściem do projektowania schematu danych jest wyjście od definicji encji i wykonywanych na nich zapytań. Należy zauważyć, że zapytanie jest w pewnym sensie nadzbiorem w stosunku do relacji. Nie tylko zawiera ono informację, że dwa obiekty są ze sobą połączone, ale reprezentuje też sposób tego połączenia pozwalający na uzyskanie wymaganych informacji. Ponieważ zapytania najczęściej odwołują się do połączonych fragmentów wielu encji podstawowym „narzędziem” modelowania Cassandra jest denormalizacja. Jest to zasadne zwłaszcza w kontekście wydajnych operacji zapisu, które oferuje baza.

Niech w przykładzie 4.6 omawiany będzie model danych dla strony internetowej, na której występują następujące odwołania do modelu danych:

1. W głównym szablonie strony przewidziane jest miejsce na nazwę użytkownika i flagę, które są wyświetlane po zalogowaniu się.

2. W panelu administracyjnym wyświetlana jest lista użytkowników, która może być filtrowana po mieście i kraju zamieszkania. Nie są na niej prezentowane pełne dane adresowe.
3. Pełny adres jest widoczny po kliknięciu na osobną podstronę profilu użytkownika.

Z punktów 1 i 2 wynika, że w kontekście użytkownika pobierane są tylko dwie składowe adresu: kraj i miasto. W przypadku punktu 3 można wykonać osobne zapytanie, gdyż adres wyświetlany jest na podstronie. Poprawny projekt modelu danych, który umożliwia ściąganie kompletu wymaganych danych w jednym zapytaniu, został schematycznie przedstawiony na diagramie 4.7.

	name	surname	addressId	city	country
982	Jakub	Turek	1842	Warsaw	Poland

Rysunek 4.7: Poprawne rozwiązanie problemu modelowania użytkownika i adresu dla zadanych wymagań.

Silna denormalizacja pozwala na budowanie efektywnych i bardzo szybkich modeli danych. W zamian wymaga ona dużego nakładu planowania. Wprowadzanie zmian do modelu zdenormalizowanego jest trudne, gdyż struktura istniejących danych jest bardzo usztywniona. W niektórych przypadkach uzysk wydajności wynikający z wykorzystania struktury zdenormalizowanej na korzyść znormalizowanej jest niewielki, a poniesione koszty sprawiają, że cały proces jest nieopłacalny. Z tego względu mechanizm OMC zapewnia szerszy wybór narzędzi do modelowania zależności pomiędzy elementami modelu danych. Narzędzia te odzwierciedlają rozważania opublikowane na blogu technologicznym portalu aukcyjnego eBay - jednego z większych użytkowników Apache Cassandra. [7].

OMC posiada trzy wbudowane tryby, które automatyzują proces zarządzania zależnościami między danymi. Różnią się one pomiędzy sobą wydaj-

nością (mierzoną w liczbie zapytań niezbędnych do pobrania kompletu wymaganych danych) oraz stopniem normalizacji:

Denormalizacja przez pole Umożliwia włączenie zdenormalizowanych pól do zamodelowanej encji. Tworzy model jednostopniowy (do pobrania kompletu informacji wymagane jest jednokrotne pobranie wiersza z tabeli).

Denormalizacja przez tabelę Umożliwia ekstrahowanie wszystkich zdenormalizowanych zależności do osobnej tabeli. Tworzy model dwustopniowy (do pobrania kompletu informacji wymagane jest dwukrotne pobranie wiersza - z tabeli nadrzędnej oraz tabeli zdenormalizowanej).

Normalizacja przez tabelę Umożliwia wiązanie zależności poprzez osobną tabelę. Tworzy model trzystopniowy (do pobrania kompletu informacji wymagane jest trzykrotne pobranie wiersza - z tabeli nadrzędnej, tabeli mapowania oraz tabeli wskazywanej).

4.4.1 Denormalizacja przez pole

Mechanizm denormalizacji przez pole pozwala wskazać encję, której wybrane pola zostaną automatycznie dołączone w wybrane miejsce modelu. Po wskazaniu encji zależnej domyślnie przenoszony jest klucz partycjonowania tej encji. Jest on również dołączany do klucza klastrowania modelu. Dzięki temu możliwe jest zdefiniowanie wielu połączeń i docieranie od obiektu nadrzędnego do obiektu zależnego.

Denormalizacja przez pole wykorzystuje typ `DenormalizedField`. Jego działanie zostanie omówione na podstawie definicji danych zawartej na diagramie 4.7. Listing 4.8 przedstawiono definicję encji adres zapisaną w OCM. Model użytkownika został pokazany na listingu 4.9.

Pole `DenormalizedField` używa dwóch parametrów konfiguracyjnych. Opcja `relates` przechowuje nazwę modelu, który jest wskazywany przez dane pole. Parametr `fields` definiuje nazwy pól encji wskazywanej, które mają

```
class Address(Model):
    address_id = UuidField(partitioning_key=True)
    street = TextField()
    street_no = IntegerField()
    flat_no = IntegerField()
    postal_code = TextField(length=32)
    city = TextField()
    country = TextField()
```

Rysunek 4.8: Denormalizacja przez pole - definicja encji adres.

```
class User(Model):
    user_id = UuidField(partitioning_key=True)
    name = TextField()
    surname = TextField()
    address = DenormalizedField(related=Address,
                                fields=['city', 'country'])
```

Rysunek 4.9: Denormalizacja przez pole - definicja encji użytkownik.

zostać włączone do definicji modelu. Ważne jest, że w celu uniknięcia kolizji nazw dołączane pola nazywane są według specjalnego klucza łączącego nazwę pola w modelu nadrzędnym z polami w encji zależnej. W podanym przykładzie dołączane pola będą miały nazwy: `address_address_id`, `address_city` oraz `address_country`.

W porównaniu do ręcznej denormalizacji użycie pola `DenormalizedField` wprowadza liczne usprawnienia:

- Cały klucz partycjonowania encji wskazywanej jest automatycznie włączany do modelu. Dzięki temu zawsze możliwe jest przejście od obiektu nadrzędnego do podrzędnego.
- Wykorzystanie zasady DRY⁶ przy definiowaniu denormalizowanych pól. Gdyby w modelu `Address` deklaracja typu pola `country` została zmie-

⁶Don't Repeat Yourself (ang. „nie powtarzaj się”) - zasada, która w kontekście programowania zaleca powielania tych samych lub bardzo podobnych fragmentów kodu w wielu miejscach.

niona na `IntegerField()` (kod kraju zamiast nazwy), typ denormalizowanego pola w modelu `User` zostałby automatycznie uaktualniony.

- Generowany jest mechanizm automatycznego rozpakowywania wartości denormalizowanych z obiektu modelu. Zakładając, że zmienna `address` przechowuje obiekt typu `Address`, a `user` obiekt typu `User`, wywołanie `user.address = address` spowoduje automatyczne przypisanie wartości pól `address_*` zmiennej `user`.
- Generowana jest metoda do automatycznego pobierania wskazywanego wiersza danych. Wywołanie `user.address.get()` zwróci obiekt typu `Address`, pod warunkiem, że klucz `user.address_address_id` jest poprawny.
- W przypadku aktualizacji danych encji zależnej mechanizm potrafi na życzenie zaktualizować również wszystkie zdenormalizowane wartości, które wskazują na ten obiekt. Można tego dokonać wywołując metodę `update(update_related=True)` na obiekcie modelu.

Pole `DenormalizedField` można używać również przy pustej liście parametru `fields`. Wtedy dowiązywany jest wyłącznie identyfikator encji zależnej. Możliwe jest również korzystanie z funkcji `get()`, która pobierze obiekt reprezentujący zależny wiersz.

Denormalizacja przez pole jest najbardziej efektywnym czasowo sposobem modelowania relacji. Wynikowy model danych jest tożsamy z przedstawionym na diagramie 4.7.

4.4.2 Denormalizacja przez tabelę

Mechanizm denormalizacji przez tabelę działa analogicznie do denormalizacji przez pola, z następującymi różnicami:

- Zdenormalizowane dane przechowywane są w osobnej tabeli, a nie w obiekcie nadrzędnym.

- Zdenormalizowane dane nie są dostępne bezpośrednio jako pola obiektu, z którego następuje odwołanie.
- Domyślnie (przy braku specyfikacji listy pól) denormalizowane są wszystkie pola obiektu zależnego.

Na listingu 4.10 przedstawiono przykład denormalizacji przez tabelę opisaną jako model OMC. Ze względu na brak wyspecyfikowanych pól (parametr `fields`) do tabeli denormalizacyjnej zostaną włączone wszystkie pola modelu `Address`. Schemat tabeli denormalizacyjnej, która powstanie na skutek zadeklarowania przykładowego modelu został zaprezentowany na diagramie 4.11.

```
class User(Model):
    user_id = UuidField(partitioning_key=True)
    name = TextField()
    surname = TextField()
    address = DenormalizedTable(related=Address, model='UserAddress')
```

Rysunek 4.10: Denormalizacja przez tabelę - definicja encji użytkownik.

	addressId	street	no	code	city	country
233	754	Nowowiejska	15/19	00-665	Warsaw	Poland

Rysunek 4.11: Denormalizacja z wykorzystaniem tabeli. Wszystkie zdenormalizowane informacje są przechowywane w osobnej tabeli.

Do obsługi zdenormalizowanej tabeli generowana jest prosta klasa modelu. Opakowuje ona dostęp do zdenormalizowanych pól. Parametr `model` specyfikuje nazwę klasy tego modelu, którą następnie można wykorzystywać w kodzie. Dostęp do wartości zdenormalizowanych dostępny jest po pobraniu obiektu zależnego. Zakładając, że zmienna `user` jest instancją klasy `User`, wywołanie `user.address.get()` zwraca dane z wiersza zdenormalizowanej tabeli. Z wyniku wykonania metody można uzyskać później wartości pól, przykładowo `user.address.get().postal_code` dla kodu pocztowego.

4.4.3 Normalizacja przez tabelę

Normalizacja przez tabelę polega na wyekstrahowaniu mapowania identyfikatorów encji wzajemnie zależnych ($id_{relates}, id_{related}$) do osobnej tabeli. Na poziomie modelu danych przypomina to wykorzystanie tabeli pośredniczącej do modelowania relacji wiele-do-wielu w relacyjnych bazach danych. Podstawową różnicą jest jednak brak zwrotności. Tabela normalizująca opisuje połączenie jednokierunkowe i nie ma możliwości efektywnego odwrócenia tego związku bez utworzenia kolejnej tabeli normalizującej w drugim kierunku.

Przykład wykorzystania normalizacji przez tabelę w modelu OMC został przedstawiony na listingu 4.12. Fizyczny model danych stworzony dla przykładowego opisu przedstawia diagram 4.13. Mapowanie powiązań w OMC przy pomocy tabeli normalizacyjnej (`user_address_norm_table`) zachodzi niejawnie. Wywołanie w kodzie aplikacji metody `user.address.get()` zwracającej listę adresów użytkownika wykonuje wiele odwołań do bazy danych. Pierwsze pobiera z tabeli normalizującej wszystkie identyfikatory powiązanych obiektów. Kolejne zapytania ściągają do pamięci szczegóły tych obiektów.

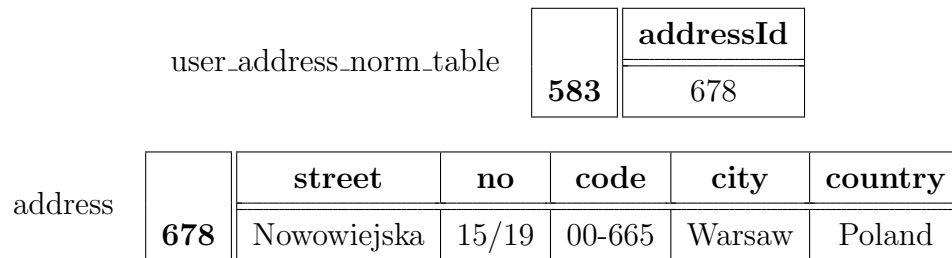
```
class User(Model):
    user_id = UuidField(partitioning_key=True)
    name = TextField()
    surname = TextField()
    address = NormalizedTable(related=Address)
```

Rysunek 4.12: Normalizacja przez tabelę - implementacja encji użytkownik.

Na podstawie opisu działania metody `get()` łatwo stwierdzić, że takie rozwiązanie bardzo źle skaluje się ze wzrostem gęstości połączeń. Przeciętna ilość zapytań niezbędnych do pobrania kompletu informacji z bazy danych wynosi:

$$|q_{avg}| = 1 + r_{avg} \quad (4.3)$$

gdzie r_{avg} jest średnią liczbą obiektów wskazywanych przez encję nadrzęd-



Rysunek 4.13: Model danych tabeli normalizacyjnej. Identyfikatorem wiersza w tabeli *user_address_norm_table* jest identyfikator użytkownika.

ną. Dla licznych, gęstych zbiorów danych oznacza to występowanie sytuacji, w których należy n osobnych zapytań dla wszystkich n elementów zapytań, które znajdują się w zbiorze.

4.4.4 Porównanie wydajności metod zarządzania zależnościami

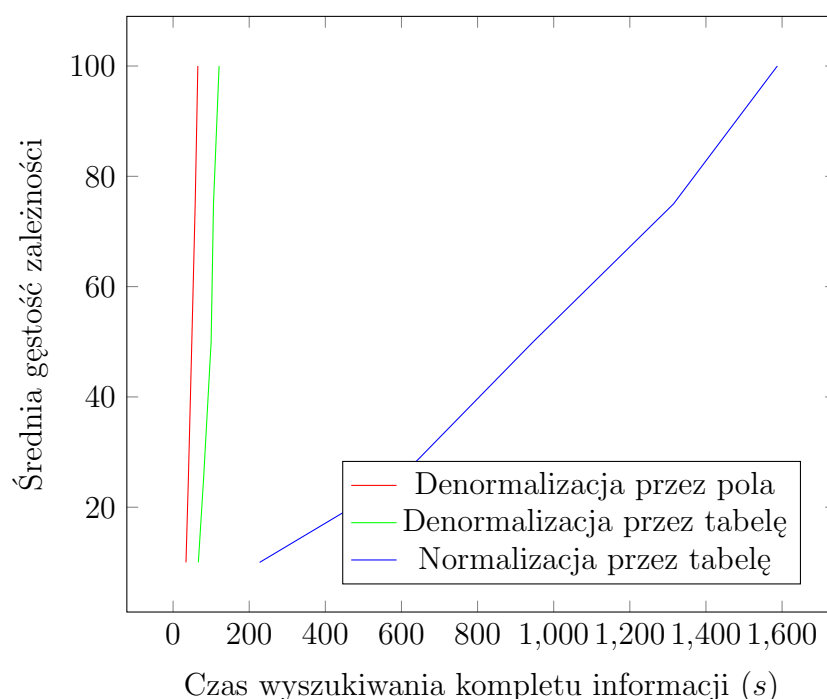
Trzy zaprezentowane metody różnią się od siebie wydajnością poszczególnych operacji w kontekście charakterystyki danych. Aby ułatwić wybór rodzaju modelowania w zależności od scenariusza użytkownika Autor przeprowadził szereg testów wydajnościowych, które ukazują zalety i wady poszczególnych typów modelowania. Wszystkie testy zostały przeprowadzone dla modelu danych użytkowników i przypisanych do nich adresów.

Skalowalność operacji wyszukiwania względem gęstości połączeń

Celem pierwszego testu jest zbadanie zależności pomiędzy wydajnością operacji pobierania kompletu danych, a gęstością połączeń pomiędzy encjami zależnymi. Gęstość połączeń to średnia ilość podzależnych elementów przypadających na jeden element nadrzędny.

W trakcie testu zostało wygenerowanych 10000 użytkowników oraz 1000 adresów. W poszczególnych iteracjach mierzone były czas pobierania wszystkich użytkowników, wraz ze wszystkimi wymaganymi informacjami adresowymi, ale dla różnej gęstości połączeń pomiędzy danymi. Wyniki testu zo-

stały przedstawione na wykresie 4.14.



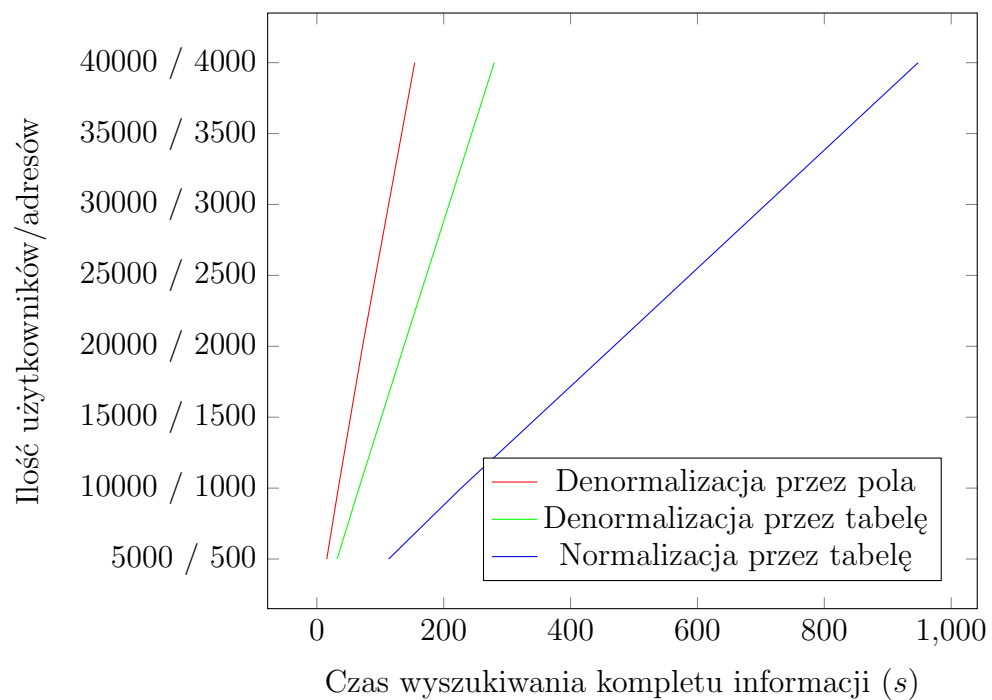
Rysunek 4.14: Porównanie czasu pobierania informacji dla różnych metod modelowania zależności (przy rosnącej gęstości połączeń).

Z wykresu można odczytać ogromną przewagę modeli zdenormalizowanych nad modelem znormalizowanym przy masowym odczycie wielu rekordów. Jest to zgodne z prognozami przedstawionymi w sekcji 4.4.3. Liczba zapytań potrzebnych do pobrania kompletu informacji z wykorzystaniem normalizacji przez tabelę rośnie wraz ze wzrostem encji w zbiorze danych, jak również z gęstością połączeń. Uzyskany wykres ma charakterystykę nieliniową. W przypadku denormalizacji przez pola/tabelę uzyskane wyniki mają charakterystykę w przybliżeniu liniową. Ilość wykonywanych zapytań zależy wyłącznie od liczby elementów w zbiorze danych.

Należy zauważyć, że wykres nie oddaje pełni różnic pomiędzy modelami. Ponieważ czas dostępu dla zależności znormalizowanej odbiega znacząco od pozostałych, skala nie odzwierciedla różnicy pomiędzy podejściami zdenormalizowanymi. Zgodnie z przewidywaniami mechanizm denormalizacji przez

pole jest dla wszystkich punktów pomiaru w przybliżeniu dwukrotnie szybszy od denormalizacji przez tabelę, ponieważ do pobrania informacji należy wykonać dwa razy mniej zapytań.

Skalowalność operacji wyszukiwania względem ilości encji Celem tego testu jest zbadanie szybkości wyszukiwania kompletu danych, kiedy gęstość połączeń w modelu jest stała, a zmienia się ilość encji nadrzędnych oraz zależnych. Wszystkie badania zostały przeprowadzone dla średniej gęstości połączeń wynoszącej 10. W kolejnych iteracjach liniowo zwiększana była ilość wierszy użytkownika i ilość adresów. Wyniki badań zostały przedstawione na wykresie 4.15.

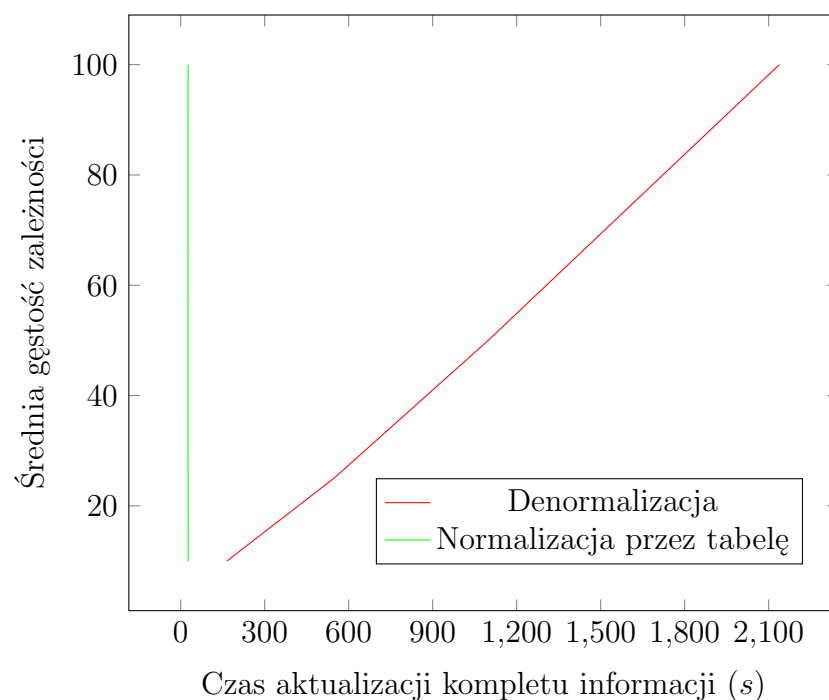


Rysunek 4.15: Porównanie czasu pobierania informacji dla różnych metod modelowania zależności (przy rosnącej ilości wierszy).

Kolejny raz efektywność zależności znormalizowanej przez tabelę znacząco odstaje od efektywności zależności zdenormalizowanych. Jednakże w przeciwieństwie do gęstości zależności, zwiększanie liczby wierszy powoduje linio-

wy wzrost czasu potrzebnego na pobranie kompletu danych. Wynika z tego, że model zdenormalizowany dużo lepiej sprawdza się w przypadku rzadkich zbiorów połączeń pomiędzy encjami. Podobnie jak w przypadku poprzedniego badania, model zdenormalizowany poprzez pola okazuje się dwa razy szybszy od denormalizacji z użyciem osobnej tabeli.

Skalowalność operacji aktualizacji względem gęstości połączeń Celem testu jest porównanie szybkości aktualizacji danych zależnych w przypadku różnych typów modelowania. W trakcie testu zostało wygenerowanych 10000 użytkowników i 1000 adresów. Zmieniana była gęstość połączeń pomiędzy encjami danych. Aktualizacji podlegały wszystkie elementy ze zbioru zależnych danych (adresy). Wyniki testu zostały przedstawione na wykresie 4.16.



Rysunek 4.16: Porównanie czasu aktualizacji informacji dla różnych metod modelowania zależności (przy rosnącej gęstości połączeń).

Z wykresu wynika, że czas aktualizacji rekordów znormalizowanych jest

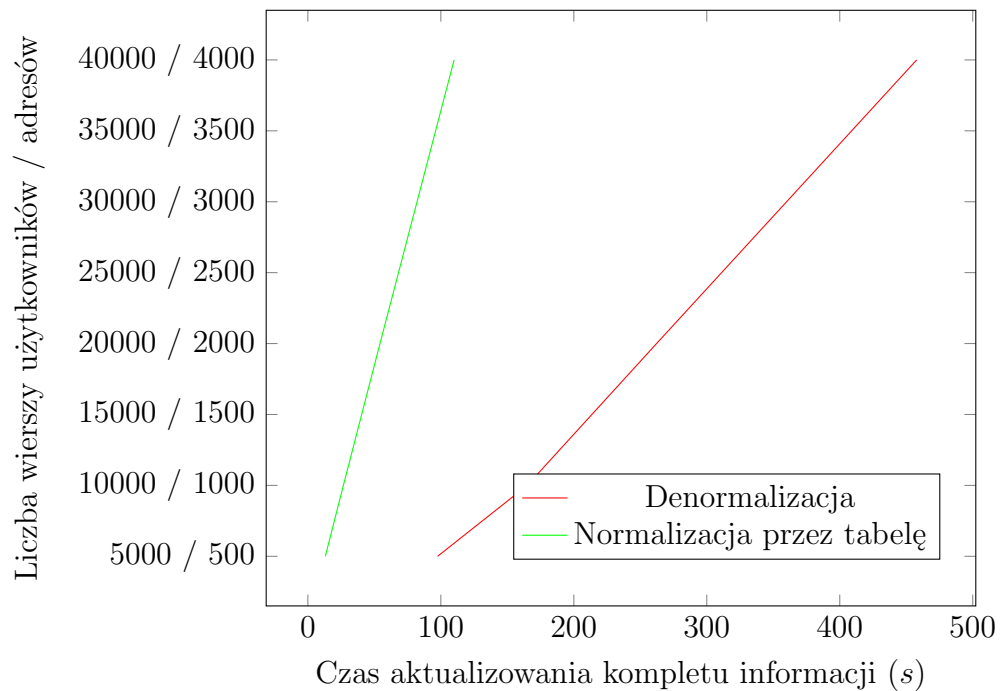
niezależny od gęstości połączeń. W modelu znormalizowanym wystarczy zaktualizować wszystkie encje zależne, więc mierzony czas faktycznie powinien być stały względem gęstości. Aktualizacja modelu znormalizowanego jest bardzo efektywna. Czas aktualizacji 1000 przedmiotów na środowisku lokalnym (przetwarzanie jednowątkowe) wynosił niewiele poniżej 30 sekund.

Oba modele zdenormalizowane mają fatalną wydajność aktualizacji kompletu danych. Problemem jest brak możliwości efektywnego wyszukiwania identyfikatorów wszystkich użytkowników, którzy są powiązani z danym adresem. Klauzula aktualizacji w Apache Cassandra nie wspiera wykorzystania indeksów drugiego stopnia. Bez wykorzystania indeksów drugiego stopnia nie jest możliwe znalezienie encji nadrzędnych, które są powiązane z daną encją zależną. Stąd wykonanie zestawu operacji `UPDATE` jest zawsze poprzedzone wyszukiwaniem `SELECT`. Kolejną przeszkodą jest ilość aktualizowanych danych. W przypadku modelu znormalizowanego wystarczyło wykonanie $|Z|$ operacji, gdzie Z to zbiór elementów zależnych (adresów). Dla modelu zdenormalizowanego liczba ta jest zwielokrotniona przez średnią gęstość połączeń pomiędzy wierszami.

Skalowalność operacji aktualizacji względem ilości encji Celem testu jest porównanie szybkości aktualizacji danych zależnych w przypadku stałej gęstości zależności dla różnych typów modelowania. Średnia gęstość w trakcie wszystkich iteracji testów wynosi 10. Liniowo zmieniają się natomiast liczby wierszy w tabelach użytkowników i adresów. Wyniki testu zostały przedstawione na wykresie 4.17.

Podobnie jak w poprzednim badaniu wydajności aktualizacji, efektywniejszy okazał się model znormalizowany. Mierzony czas rośnie proporcjonalnie do ilości nadpisywanych rekordów. Dla zadanego problemu jest to rozwiązanie optymalne.

W przypadku denormalizacji można zauważyć, że wzrost liczby wierszy w tabeli ma łagodniejszy wpływ na charakterystykę wykresu czasu aktualizacji niż rosnąca gęstość połączeń. Przykładowo dwukrotny wzrost gęstości



Rysunek 4.17: Porównanie czasu aktualizacji informacji dla różnych metod modelowania zależności (przy rosnącej ilości wierszy).

połączeń odnosi podobny skutek czasowy co ośmiokrotny wzrost ilości danych.

4.4.5 Wnioski z analizy wydajności modelowania zależności

Nie ma uniwersalnego sposobu, który pozwala na uzyskanie optymalnego wyniku przy mapowaniu zależności w Apache Cassandra. Dobór metody powinien zależeć od charakterystyki przechowywanych danych.

Szukając pewnego rodzaju uogólnienia Autor wprowadził pomocniczą definicję typowego zbioru danych. Posiada on zbiór następujących cech:

- Duża liczba wierszy w zbiorze (> 100000).
- Operacje wyszukiwania są znacznie częstsze niż aktualizacji. Stosunek

częstotliwości ich występowania wynosi $20 \div 1$ (istnieje 5% szansy, że użytkownik przeglądający profil zaktualizuje dane adresowe).

- Gęstość występowania zależności jest niewielka. Na jeden wiersz encji nadrzędnej przypada średnio nie więcej niż dziesięć wierszy encji zależnej.

Dla tak określonego zbioru danych zdecydowanie najlepszym wyborem jest modelowanie z wykorzystaniem denormalizacji poprzez pola. Głównym kryterium jest tutaj znacząca przewaga w czasie wyszukiwania danych i akceptowalny narzut przeznaczony na aktualizację danych. Liczne artykuły prezentujące studium przypadku użycia Apache Cassandra każą sądzić, że ten rodzaj modelowania zależności jest prawdopodobnie najlepiej dostosowany do specyfiki schematu danych.

W przypadku małych zbiorów danych (preferowana ilość wierszy w tabelach ≤ 10000) o niskiej gęstości zależności warto rozważyć użycie reprezentacji znormalizowanej. Nie jest ona obciążona narzutem na aktualizację danych. Zaletą tej reprezentacji jest większa elastyczność schematu. W przeciwieństwie do postaci zdenormalizowanej, gdzie może zdarzyć się, że przy dodaniu kolejnego zapytania trzeba rozszerzyć tabelę denormalizacji o kolejną kolumnę, zapewnia on dostęp do kompletu danych bez konieczności wprowadzania zmian w strukturze.

Reprezentacja zdenormalizowana jest zalecana w przypadku, gdy wyszukiwanie danych nie jest kluczowym kryterium. Przykładem mogą być dane służące do generacji raportów, gdzie liczba aktualizacji może znacznie przewyższać liczbę wyszukiwań.

Duży problem stanowi modelowanie zależności o dużej liczbie połączeń między danymi. Wraz ze wzrostem gęstości połączeń możliwości zawężają się do:

- Zaniedbania aktualizacji zdenormalizowanych danych. W praktyce jest to bardzo kłopotliwe i może być zastosowane jeżeli udaje się zdenormalizowanie wyłącznie danych przeznaczonych tylko do odczytu.

- Poświęcenia czasu aktualizacji kosztem czasu dostępu do danych i wyboru modelu zdenormalizowanego.
- Poświęcenia czasu dostępu do danych kosztem czasu aktualizacji i wyboru modelu znormalizowanego.
- Zmiany modelowania danych. Przekształcony schemat może między innymi zredukować zagęszczenie zależności lub konieczność wykonywania operacji aktualizacji.

Jeżeli żadna z powyższych możliwości nie daje się zastosować lub nie wprowadza akceptowalnej poprawy wydajności należy rozważyć wybór innego silnika bazodanowego. Dla zbiorów o bardzo gęstej siatce zależności naturalne zastosowanie mają grafowe bazy danych. Przykładem dojrzałego projektu silnika grafowego o dobrych własnościach skalowania jest Neo4j⁷.

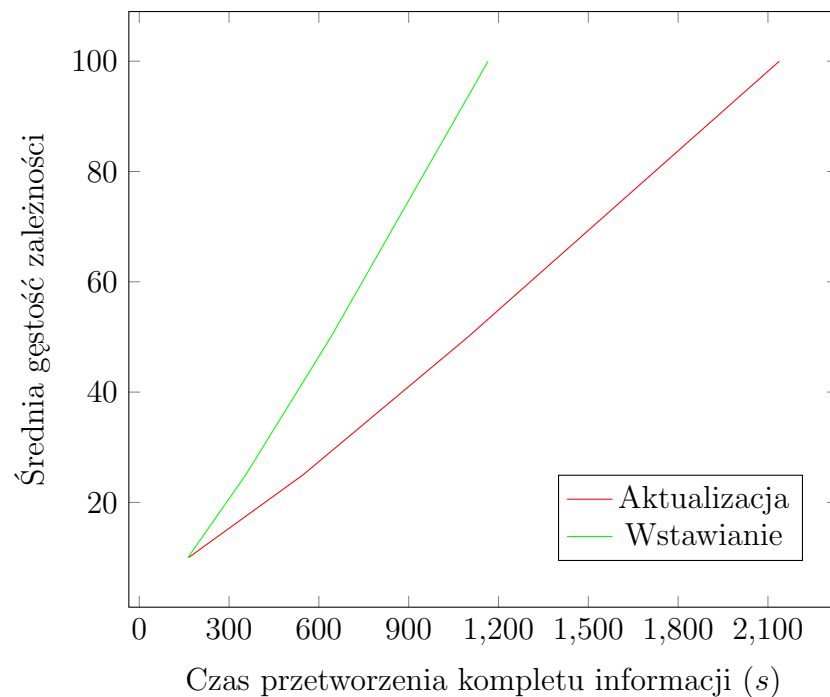
W omawianym w sekcji 4.4 przykładzie można wyeliminować operacje aktualizacji danych przy specyficznych założeniach:

- Żadna encja nie jest powiązana z użytkownikiem jako encja nadrzędna.
- W trakcie aktualizacji adresu aplikacja dysponuje pełnym kompletem informacji o użytkowniku.
- Nie istnieje ryzyko konfliktów synchronizacji (nie jest możliwa równoczesna edycja tego samego profilu przez różne instancje aplikacji).
- Zdublikowanie istniejącego adresu nie ma wpływu na poprawność danych.

Jeżeli powyższe założenia są spełnione aktualizacja może być zastąpiona poprzez wstawienie nowych danych do bazy. Ta operacja jest bardziej wydajna ze względu na pominięcie etapu wyszukiwania użytkowników, którzy są przypisani do danego adresu. Porównanie szybkości aktualizacji i wstawiania

⁷Strona domowa projektu Neo4j dostępna jest pod adresem <http://www.neo4j.org>.

nowych danych do bazy (z unieważnianiem poprzednich) zostało przedstawione na wykresie 4.18. Testy zostały przeprowadzone na środowisku lokalnym dla stałej liczby wierszy w tabelach (10000 użytkowników oraz 1000 adresów). W kolejnych iteracjach modyfikowana była średnia gęstość zależności w modelu. Modyfikacji / wstawianiu podlegały wszystkie wiersze.



Rysunek 4.18: Porównanie czasu wstawiania i aktualizacji informacji (przy rosnącej gęstości połączeń).

Wyniki pokazują, że dla gęstych modeli danych operacja wstawiania może być nawet dwukrotnie wydajniejsza niż aktualizacja. Zastosowanie wielowątkowej aplikacji do przeprowadzenia testów mogło dalej pogłębić tę różnicę. Z wyniku testu można wywnioskować, że odnajdywanie/nakładanie odpowiednich ograniczeń na dane i przekształcanie modelu może znacząco poprawiać wydajność również w przypadkach, w których w ogólności Cassandra nie sprawdza się zbyt dobrze.

Na podstawie porównania wydajności różnych metod modelowania zależności między danymi można stwierdzić, że nie ma przypadku, w którym de-

normalizacja z wykorzystaniem tabeli ma przewagę nad pozostałymi dwoma metodami. W każdym przypadku może być ona zastąpiona denormalizacją przez pola/modeliem znormalizowanym z uzyskaniem identycznej lub lepszej wydajności. Stosowanie tego rodzaju relacji jest uzasadnione jeżeli:

- Uzyskana wydajność jest zadowalająca, a istotne jest fizyczne uporządkowanie i wydzielenie danych (na przykład gdy część systemu wykorzystuje wydruki w formie surowej).
- Opisywane są bardzo szerokie wiersze i problemem jest przekroczenie dwóch miliardów kolumn⁸.

4.5 Wsparcie dla wzorców modelowania

Jednym z pryncypiów projektowych mechanizmu OMC jest obsługa wzorców modelowania. W literaturze zostały zebrane i opisane liczne przykłady efektywnych i nieefektywnych projektów modelu danych dla Apache Cassandra. Z przeanalizowanych opisów Autor wybrał zbiór reużywalnych wzorców. Następnie, dla każdego wzorca, przygotowany został mechanizm, który ułatwia jego modelowanie w środowisku OMC. Dzięki temu mechanizm został wzbogacony o obsługę:

- Szeregów zdarzeń, czyli zbioru uporządkowanych względem czasu wiadomości.
- Kolejek *First In, First Out*, *Last In, Last Out* i podobnych.
- Indeksów wartości unikalnych.
- Aktualizacji selektywnej.

⁸Apache Cassandra narzuca limit dwóch miliardów kolumn w jednym wierszu.

4.5.1 Szereg zdarzeń

Apache Cassandra została zaprojektowana z myślą o wspieraniu szeregów zdarzeń. Szereg zdarzeń to uporządkowane względem czasu wystąpienia wpisy do bazy danych, które opcjonalnie mogą posiadać z góry określony czas życia. Typowym przykładem szeregu zdarzeń są wartości odczytów z urządzeń pomiarowych. Wartości te zwyczajowo są pobierane cyklicznie, a do wykorzystania w czasie rzeczywistym przydatne są tylko najnowsze odczyty. Doskonały opis modelowania szeregów zdarzeń przedstawiony jest we wpisie Patricka McFadina na portalu **Planet Cassandra**. [13] Artykuł analizuje przykład stacji meteorologicznych. Przedstawia trzy sposoby zapisu szeregu zdarzeń w Apache Cassandra.

Odczyty dla tego samego urządzenia w pojedynczym wierszu Najprostszą możliwością jest zebranie wszystkich próbek dla danego urządzenia w jednym wierszu tabeli. Kolejne wartości zapisywane są ze śladem czasu w nazwie kolumny. Listing 4.19 przedstawia wzorzec zapisany jako model OMC.

```
class TimeSeriesPatternOne(Model):  
    weatherstation_id = TextField(partitioning_key=True)  
    event_time = TimestampField(clustering_key=True, auto_on_create=True)  
    temperature = TextField()
```

Rysunek 4.19: Modelowanie szeregu zdarzeń. Zapis kolejnych odczytów jako kolumny.

Zapis śladu czasu w nazwie kolumny można uzyskać poprzez dołączenie pola do złożonego klucza głównego tabeli jako *clustering key*⁹. Dołączenie to dokonuje się poprzez ustawienie flagi `clustering_key` na `True`.

Udogodnieniem dla programisty korzystającego z OMC jest możliwość automatycznej generacji śladu czasu dla chwili obecnej. Ślad czasu może

⁹Clustering key (ang. dosłownie „klucz klastrowania”) - część złożonego klucza głównego tabeli, która określa w jakiej kolejności grupowane są wpisy względem klucza partycjonowania (identyfikatora wiersza).

być wygenerowany w trakcie tworzenia obiektu (`auto_on_create=True`) lub w trakcie wstawiania obiektu do bazy danych (`auto_on_save=True`).

Grupowanie odczytów dla urządzenia i fragmentu daty W przypadku gdy odczytów jest zbyt dużo by mieściły się w jednym wierszu sugerowanym rozwiązaniem jest dodatkowe partycjonowanie wierszy po komponencie daty: miesiącu, dniu, itd. Autor wzorca sugeruje utworzenie w modelu dodatkowego pola, do którego wpisywany jest komponent daty i dołączenie tego pola do klucza partycjonowania. Takie modelowanie dziedziny może prowadzić do potencjalnych błędów. Ta sama dana - ślad czasu - przechowywana jest w dwóch polach, które różnią się formatowaniem.

OMC posiada mechanizm partycjonowania dat po komponencie. W modelu znajduje się tylko jedna dana - faktyczny ślad czasu. OMC na jej podstawie automatycznie tworzy i formatuje klucz partycjonowania. Dzięki temu programista nie musi dbać o ręczne wyznaczanie komponentu daty. Listing 4.20 przedstawia opis modelu z partycjonowaniem względem dnia.

```
class TimeSeriesPatternTwo(Model):
    weatherstation_id = TextField(partitioning_key=True)
    event_time = TimestampField(clustering_key=True,
                               auto_on_create=True,
                               partitioning_by_day=True)
```

Rysunek 4.20: Modelowanie szeregu zdarzeń. Partycjonowanie zdarzeń po komponencie daty.

Pole `event_time` modelu `TimeSeriesPatternTwo` jest kluczem klastrowania. Dodatkowo posiada ono ustawioną flagę `partitioning_by_day=True`, która spełnia trzy zadania:

- Przy tworzeniu schematu tabeli dodaje automatycznie nową kolumnę `event_time_day`¹⁰ o typie tekstowym. Pole to jest automatycznie włączone do klucza partycjonowania.

¹⁰Nazwa jest tworzona z nazwy pola oraz nazwy komponentu daty.

- Przy zapisywaniu obiektu do bazy danych automatycznie wypełnia wartość dodatkowego pola o ślad czasu obciążenia do z dokładnością do dnia.
- Przy wyszukiwaniu obiektów po dacie automatycznie filtruje dane zarówno po wartości pola `event_time_day`, jak i `event_time`.

Dostępne jest kilka rodzajów komponentów daty, po których można dokonać partycjonowania. Różnią się one ostatnim członem nazwy parametru. Wybór obejmuje partycjonowanie po latach (`_year`), miesiącach (`_month`), dniach (`_day`), godzinach (`_hour`), minutach (`_minute`), a także sekundach (`_second`).

Odczyty z ograniczoną pamięcią Korzystając z wbudowanego mechanizmu Cassandra możliwe jest tworzenie automatycznie przedawniających się wartości. Wstawiając pozycję do bazy danych można ustawić czas życia, po którym zostanie ona usunięta. Przedawnianie się wartości może być wykorzystywane w połączeniu z odwrotną kolejnością sortowania wpisów (od najnowszego do najstarszego). Dzięki temu bardzo łatwo jest budować widoki, które prezentują kilka najnowszych informacji.

OMC zapewnia dwa rodzaje wsparcia dla czasu życia wartości. Czas życia może być ustawiany dla indywidualnych obiektów. Przy zapisywaniu obiektu należy na jego instancji wywołać metodę `save(ttl=20)`. Parametr `ttl` definiuje czas życia wartości w sekundach. Czas życia można również ustawić globalnie dla wszystkich instancji modelu. Wykorzystuje się do tego pole `__ttl__` klasy. Zaprezentowano to na listingu 4.21.

Dla każdego pola, które jest częścią klucza klastrowania można ustawić odwrotną kolejność sortowania elementów. Do tego celu należy ustawić flagę `descending_clustering` na `True`.

```
class TimeSeriesPatternThree(Model):
    __ttl__ = 20
    weatherstation_id = TextField(partitioning_key=True)
    event_time = TimestampField(clustering_key=True,
                                descending_clustering=True,
                                auto_on_create=True)
    temperature = TextField()
```

Rysunek 4.21: Modelowanie szeregu zdarzeń. Przedawniające się odczyty jako kolumny z sortowaniem od najnowszego.

4.5.2 Wsparcie dla kolejek

W sekcji 2.8 został przytoczony antywzorzec kolejki. Problem został opisany przez Alekseya Yeschenko na blogu firmy DataStax. [9] Wynika on ze sposobu obsługi operacji usuwania przez silnik Apache Cassandra. Usuwanie nie jest realizowane natychmiast. Zamiast tego kasowana kolumna oznaczana jest specjalnym znacznikiem *tombstone* i rezyduje w pamięci aż do momentu, gdy nastąpi właściwe wyrzucanie. Przy dużej liczbie usunięć kolumn z wiersza może zdarzyć się, że zapytanie filtrujące z limitem będzie działało znacznie wolniej niż dla wiersza z identyczną ilością kolumn, ale bez poprzedzających usunięć. Wynika to z faktu, że stwierdzenie, że dana wartość jest „martwa” wymaga pobrania kolumny i zdeserializowania jej wartości. Jeżeli po deserializacji wykryty zostanie znacznik *tombstone*, wtedy pobierana jest kolejna wartość. Proces powtarzany jest aż do skutku.

Kolejka może być zamodelowana efektywniej w przypadku, gdy możliwe jest wskazanie miejsca, w którym rozpoczynają się nieusunięte kolumny. Przykładem może być kolejka typu FIFO¹¹. Wyszukując kolejny element w kolejce można nałożyć warunek, że czas wstawienia elementu do kolejki jest późniejszy niż czas ostatnio usuniętego elementu. OMC dostarcza wsparcia dla takich kolejek, co zostało przedstawione na listingu 4.22.

Parametr `__track_deletes__` umożliwia śledzenie operacji usuwania.

¹¹FIFO (First In, First Out) - rodzaj kolejki, w którym elementy obsługiwane są według kolejności wstawienia.

```
class FIFOQueue(Model):
    __track_deletes__ = ('enqueued_at', 'ASC')
    name = TextField(partitioning_key=True)
    enqueued_at = UuidField(type=TimeUuid,
                           auto_generate=True,
                           clustering_key=True)

    payload = DataField()

    # create three queue elements
    first_element = FIFOQueue(name = '1', payload = 'firstPayload')
    second_element = FIFOQueue(name = '2', payload = 'secondPayload')
    third_element = FIFOQueue(name = '3', payload = 'thirdPayload')
    # save elements into database
    first_element.save()
    second_element.save()
    third_element.save()
    # delete first two elements from database
    first_element.delete()
    second_element.delete()
    # get element with name '3'
    FIFOQueue.objects().find(name='3')[:1]
```

Rysunek 4.22: Modelowanie kolejki typu FIFO w OMC.

Przyjmuje on parę łańcuchów znaków: nazwę pola, którego wartość jest śledzona podczas usuwania oraz warunek zapisywania śledzonej wartości. Na chwilę obecną obsługiwane jest śledzenie pól typu `UuidField` (w trybie `timeuuid`) oraz `TimestampField`, a także dwa warunki śledzenia: `ASC` oraz `DESC`. Warunek *ascending* oznacza, że śledzona wartość jest aktualizowana tylko w przypadku, gdy jest mniejsza od usuwanej obecnie wartości. Warunek *descending* działa odwrotnie - aktualizuje śledzoną wartość tylko w przypadku, gdy jest większa od usuwanej obecnie.

Ustawienie parametru `__track_deletes__` powoduje, że śledzona wartość jest automatycznie dodawana do wszystkich zapytań wyszukiwanych obiekty. Oznacza to, że ostatnie wywołanie w listingu 4.22 będzie miało ustawione trzy warunki filtrowania:

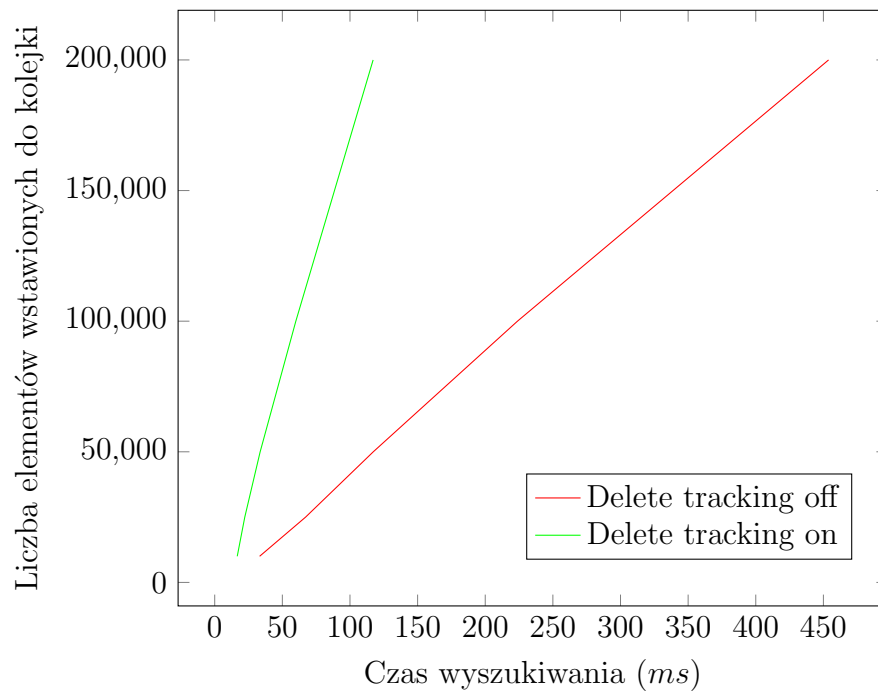
- Nazwa elementu (`name`) jest równa `'3'` - wynika z `find(name='3')`.
- Wybierz maksymalnie jeden element (`LIMIT 1`) - wynika z zakresu `[:1]`.
- Czas dodania do kolejki (`enqueued_at`) jest większy bądź równy wartości `second_element.enqueued_at` - automatycznie dodawane przez OMC.

Wykorzystanie opcji `__track_deletes__` w sposób znaczący wpływa na wydajność kolejki FIFO. Na wykresie 4.23 przedstawiono porównanie szybkości działania wyszukiwania pierwszego nieusuniętego elementu kolejki przy włączonym oraz wyłączonym śledzeniu. Pomiar wydajności został przeprowadzony na lokalnej konfiguracji klastra Apache Cassandra. Czas mierzony był przez aplikację wykorzystującą mechanizm OMC. Aplikacja ta wstawiała do kolejki y przedmiotów, po czym usuwała pierwsze $y - 1$ przedmiotów z tej kolejki. Mierzony był czas wyszukiwania pierwszego elementu należącego do kolejki (z ograniczeniem `LIMIT 1`, podobnie jak przy wywołaniu metody `find` w listingu 4.22).

Pomimo znaczącej poprawy czasu dostępu do danych przy wykorzystaniu parametru `__track_deletes__` można zauważyć, że modelowanie nie jest optymalne. W przypadku braku usunięć danych z tabeli czas dostępu do pojedynczej danej powinien być w przybliżeniu stały. Z wykresu wynika jednak, że wyszukiwanie wydłuża się wraz ze wzrostem liczby elementów w kolejce. Pozwala to wyciągnąć wniosek, że efektywne modelowanie w Cassandrze powinno całkowicie eliminować lub radykalnie ograniczać ilość wykonywanych operacji usunięć.

4.5.3 Selektywna aktualizacja

Apache Cassandra jest przystosowana do obsługi wielu równoległych źródeł danych. Ze względu na brak transakcyjności jednym ze sposobów na uniknięcie problemów synchronizacji jest zalecana przez Dave'a Gardnera selektywna



Rysunek 4.23: Porównanie czasu wstawiania wielu rekordów.

aktualizacja. [10] Edytując zawartość danego wiersza należy przesyłać do bazy danych wyłącznie wartości, które rzeczywiście uległy modyfikacji. Inaczej istnieje prawdopodobieństwo niecelowego nadpisania zmian wprowadzonych równolegle przez inne źródło danych. OMC domyślnie korzysta z mechanizmu selektywnej aktualizacji, co przedstawiono na listingu 4.24.

W pierwszej kolejności metoda `save()` sprawdza, którą z dwóch operacji, aktualizacji czy wstawienia obiektu, należy wykonać. Przy wykonywaniu aktualizacji do bazy danych przesyłana jest tylko wartość zmienionych pól, w danym przykładzie tylko priorytetu. Jeżeli pomiędzy pobraniem obiektu, a aktualizacją priorytetu i zapisaniem obiektu w bazie danych inny węzeł dokona zmiany opisu obiektu, opis ten nie zostanie ponownie nadpisany poprzednią wartością.

Mechanizm aktualizacji selektywnej można wyłączyć. Można dokonać tego wywołując metodę `save()` z parametrem `selective_update=False`. Alternatywnie selektywna aktualizacja może zostać wyłączona dla całego mo-


```

class SelectiveUpdate(Model):
    name = TextField(partitioning_key=True)
    description = TextField(length=1024)
    priority = IntegerField()
    # find object with name 'su_test'
    su_test = SelectiveUpdate.objects().find(name='su_test').get()
    # update priority
    su_test.priority = 1
    # update object
    su_test.save()

```

Rysunek 4.24: Przykład selektywnej aktualizacji.

delu poprzez zadeklarowanie pola `__selective_update__ = False`.

4.5.4 Indeks wartości unikalnych

Aby umożliwić wyszukiwanie po wartościach kolumny, która nie należy do klucza głównego tabeli należy użyć indeksu drugiego rzędu (*secondary index*). Według dokumentacji Apache Cassandra [14] indeksy drugiego stopnia powinny być używane na kolumnach, które silnie grupują wiersze encji nadrzędnej. Wynika to z wewnętrznej implementacji indeksów drugiego stopnia. Załóżmy, że w systemie istnieje tabela przechowująca dane użytkowników, w której zapisane są wiersze zaprezentowane na diagramie 4.25.

	name	surname	email	city
jturek	Jakub	Turek	J.Turek@pw.edu.pl	Warsaw
	name	surname	email	city
manisero	Michał	Aniserowicz	M.Aniserowicz@pw.edu.pl	Warsaw

Rysunek 4.25: Przykładowe wartości w tabeli użytkowników.

Przykładowy wygląd danych w tabeli, która realizuje indeks drugiego stopnia dla miasta (kolumna *city*), został zaprezentowany na diagramie 4.26.

Wadą rozwiązania przedstawionego na diagramie 4.26 jest brak skalowalności. Przy dodawaniu użytkowników wiersz rozszerza się o kolejne kolumny,

	jturek	manisero
Warsaw	null	null

Rysunek 4.26: Przykładowa tabela, która mogłaby realizować indeks drugiego stopnia dla kolumny *city* z diagramu 4.25.

które przechowywane są na tym samym węźle. Zajmuje to dużo miejsca na jednej partycji. Dodatkowo wszystkie zapytania o dany indeks zawsze odwołują się do jednego węzła, co skutkuje nierównomiernym rozłożeniem obciążenia. Ze względu na wymienione wady indeksy drugiego stopnia są realizowane lokalnie. Każdy węzeł zawiera informacje o wszystkich indeksowanych kolumnach, ale wyłącznie o kluczach wierszy, które znajdują się na danym węźle. Oznacza to, że jeżeli użytkownicy **jturek** oraz **manisero** znajdują się na różnych węzłach, również indeks jest równomiernie dystrybuowany pomiędzy te dwa węzły.

Lokalne indeksowanie wymusza odwołania do wszystkich węzłów dla każdego wyszukiwania. Zakładając, że indeksowanych wierszy (zarówno wszystkich, jak i tych z wyniku zapytania) jest dużo więcej niż węzłów, nadal znacząco przyspiesza to wyszukiwanie. Sytuacja zmienia się, gdy indeksowana jest wartość unikalna. Wtedy wbudowany w Cassandrę mechanizm wymusza wykonanie wielu zapytań zamiast jednego.

Indeks wartości unikalnych można zamodelować samodzielnie w prosty sposób. [15] Wystarczy utworzyć nową tabelę, w której unikalna wartość jest identyfikatorem. W wierszu dla danej unikalnej wartości zapisany jest odpowiedni identyfikator. Przykład takiego indeksu dla adresu e-mail (kolumna *email*) użytkowników przedstawiono na diagramie 4.27.

OMC wspiera automatyczne tworzenie indeksów dla wartości unikalnych. Na listingu 4.28 przedstawiono sposób utworzenia indeksu dla wartości unikalnych dla pola *email*.

Flaga `searchable_unique=True` powoduje, że dla danego pola tworzony jest indeks unikalny. Indeks jest w całości zarządzany przez OMC. Podczas wywołania `User.objects().find(email='J.Turek@pw.edu.pl')` me-

UserEmailIndex	J.Turek@stud.elka.pw.edu.pl	jturek
		null
	M.Aniserowicz@stud.elka.pw.edu.pl	manisero
		null

Rysunek 4.27: Przykład tabeli przechowującej indeks dla unikalnych wartości (adres e-mail użytkownika).

```
class User(Model):
    id = UuidField(type=TimeUuid, auto_generate=True, partition_key=True)
    name = TextField()
    surname = TextField()
    email = TextField(searchable_unique=True)
    city = TextField(searchable=True)
```

Rysunek 4.28: Modelowanie tabeli użytkownicy z wykorzystaniem indeksu dla wartości unikalnych.

chanizm wyszukuje identyfikator użytkownika w stworzonym przez siebie indeksie. OMC aktualizuje indeksy podczas wykonywania operacji `save()` oraz `delete()`.

Zarówno indeks wartości unikatowych, jak również domyślny indeks drugiego stopnia (tworzony przy użyciu flagi `searchable=True`) nazywane są przy użyciu tej samej konwencji. Pierwszą składową jest nazwa encji, drugą - nazwa pola, a trzecia to słowo „index”. Dla przykładu 4.28 indeks wartości unikatowych zostanie umieszczony w tabeli o nazwie `user_email_index`.

4.6 Przetwarzanie partiami

Apache Cassandra zapewnia wsparcie dla operacji atomowych. [16] Wykorzystując partie¹² można wykonać zestaw operacji, z których (cytując dokumentację) „jeżeli jedna się powiedzie, to wykonają się wszystkie”. Nie jest to jednak odpowiednik transakcji znanych z relacyjnych baz danych. Apache

¹²Partia (ang. batch) - zgrupowany zestaw zapytań INSERT, UPDATE lub DELETE.

Cassandra nie zapewnia integralności wykonywanych w partii zapytań: dane mogą być modyfikowane równolegle przez inne zapytania. Przykład atomowych operacji został zaprezentowany na diagramie 4.29. Jest to zestaw dwóch zapytań, które usuwają z bazy danych użytkownika i właściwy mu indeks dla pola e-mail.

```
BEGIN BATCH
DELETE FROM user WHERE userId = 'jturek'
DELETE FROM user_mail_index WHERE mail = 'J.Turek@pw.edu.pl'
APPLY BATCH;
```

Rysunek 4.29: Przykład operacji na partii, która usuwa z bazy danych użytkownika *jturek* i jego indeks dla pola email.

Silnik wykorzystywany w OCM zapewnia wsparcie dla przetwarzania partiami. Potrafi kolejkować wszystkie zapytania DML¹³, a następnie wykonywać zapytania z kolejki jako operacja atomowa. Przykład wykorzystania operacji na partiach w OCM zaprezentowano na listingu 4.30.

```
# starts batch
Model.begin_batch()
# create new Item
item = Item(name='item', description='desc', price=12.5)
item.save()    # save is queued
wishlist = Wishlist(userId='jturek', itemId=item.id)
wishlist.save()    # save is queued
# executes item.save() and wishlist.save()
Model.apply_batch()
```

Rysunek 4.30: Przykład wykonania operacji na partii, która wstawia do bazy nowy przedmiot i dodaje go do listy życzeń użytkownika.

Wywołanie metody `Model.begin_batch()` powoduje rozpoczęcie kolejkowania zapytań DML. Operacja wstawienia nowego przedmiotu (`item.save()`)

¹³Data Modification Language (ang. język modyfikacji danych) - podzbiór zapytań obejmujących wstawianie (INSERT), aktualizację (UPDATE) oraz usuwanie (DELETE) danych.

nie wykonuje się od razu, ale jest umieszczana w kolejce. Do kolejki trafia również operacja aktualizacji listy życzeń użytkownika (`wishlist.save()`). Wywołanie metody `Model.apply_batch()` powoduje, że partia dwóch zgrupowanych zapytań wykonuje się.

4.7 Wsparcie dla liczników

Apache Cassandra udostępnia wsparcie dla specjalnego typu danych - liczników. [17] Liczniki to kolumny specjalnego typu `counter`, na których można wykonywać tylko dwie operacje: inkrementacja lub dekrementacja. Poza specjalnym typem zmienia się również sposób interakcji z tabelą. Nie można wstawiać do niej wartości, a jedynie aktualizować liczniki dla różnych kombinacji pozostałych pól.

OMC wspiera modelowanie liczników. Przykład przedstawiono na listingu 4.31.

```
class TripCounter(Model):
    country = TextField(partitioning_key=True)
    visits = CounterField()
    sweden = TripCounter(country='Sweden')
    sweden.visits.increment(1)
    sweden.save()
    poland = TripCounter(country='Poland')
    poland.visits.increment(5)
    poland.visits.decrement(2)
    poland.save()
    s_visits = TripCounter.find(country='Sweden').get().visits      # 1
    p_visits = TripCounter.find(country='Poland').get().visits      # 3
```

Rysunek 4.31: Przykład użycia liczników w OCM.

Obsługa modelu posiadającego licznik różni się od wykorzystania standardowego modelu. Przede wszystkim wartości wszystkich pól są ustawione w trybie `read only`, czyli nie jest możliwe wykonanie operacji podstawienia `sweden.country = 'Germany'`. Dzięki temu wszystkie operacje wyko-

nywane na liczniku są zawsze odpowiednio interpretowane, gdyż nie może się zmienić ich kontekst. W przeciwnym razie trudno byłoby ocenić intencje programisty, który najpierw inkrementuje licznik o 2, a następnie zmienia wartość kraju dla obiektu ze zwiększonym licznikiem. Mogło być to działanie celowe (operacje wykonane w nieintuicyjnej kolejności) lub zwykły błąd.

Pole typu `CounterField` udostępnia dwie operacje: inkrementację wartości licznika (`increment()`) oraz dekrementację (`decrement()`). Metody wywołane bez podania parametru standardowo zwiększają/zmniejszają wartość o 1. Obie funkcje są łączne. Oznacza to, że tak jak w przykładzie można dokonywać wielu operacji na jednym liczniku, które przed uaktualnieniem obiektu w bazie danych są sumowane.

Po wykryciu pola licznikowego w definicji modelu OCM automatycznie zmienia sposób obsługi obiektu. Nowo powstały obiekt jest traktowany tak, jakby został wcześniej pobrany z bazy danych - operacja `save()` dokonuje aktualizacji, a nie wstawienia.

4.8 Migracje pomiędzy wersjami modelu

Typowym problemem związanym z wykorzystaniem baz danych jest aktualizacja schematu. Dotyczy on również narzędzi do mapowania obiektowego. Problem ten pojawia się, gdy udało się wdrożyć pewien schemat i istnieją już dane na nim operujące, a wymagana jest modyfikacja modelu. Typowe scenariusze obejmują między innymi dodawanie kolumn do istniejących tabel, dodawanie nowych modeli oraz usuwanie istniejących tabel.

Modyfikacje wprowadzane do działającego modelu zazwyczaj nie są skomplikowane i łatwo wprowadzić je ręcznie. Dla przedstawionej propozycji mapowania obiektowego dla Cassandra można jednak znaleźć przypadki, których automatyzacja zaoszczędzi dużo pracy. Niech w modelu danych będzie wykorzystane pole typu `DenormalizedField`. Jeżeli wymagane będzie dodanie nowej zdenormalizanej danej należy:

1. Rozszerzyć model danych o nową kolumnę.

2. Dla każdego połączenia pobrać daną źródłową i przepisać ją do zdenormalizowanego pola.

Dzięki wbudowanemu wsparciu dla migracji OMC pozwala wykonywać takie operacje automatycznie. Dedykowany mechanizm pozwala tworzyć punkty migracyjne, czyli kolejne wersje modelu danych. Potrafi również odczytywać wersję modelu na podstawie schematu bazy danych i aktualizować ją do najnowszej (lub wybranej) rewizji.

Mechanizm migracji w OMC jest oparty na najnowszej wersji projektu Django. [23] Podstawową jednostką migracji jest klasa `Migration`. Definiuje ona dwa pola:

zależności czyli skrypty migracyjne, które muszą zostać wykonane przed uruchomieniem operacji,

operacje czyli zmiany w modelu, które należy wykonać aby zaktualizować model danych.

Niech dany będzie model użytkownika przedstawiony na listingu 4.32. W pierwszej wersji modelu identyfikator jest typu `text`. Poza tym użytkownik posiada jeszcze imię typu tekstowego. W drugiej wersji modelu dodano pole tekstowe nazwisko oraz zmieniono typ identyfikatora na `timeuuid`.

Pierwszym krokiem jest stworzenie początkowej migracji. Odpowiada za to metoda `Model.migrations().make(initial=True)`. Wynikiem będzie powstanie nowego pliku `migration_initial.py`. Zakładając, że polecenie zostało wywołane dla pierwszej wersji modelu, zawartość tego pliku jest przedstawiona na listingu 4.33.

Po utworzeniu migracji początkowej i zmianie modelu danych należy utworzyć kolejny punkt migracji. Wywołanie metody `make()` utworzy plik `migration_0001.py`, którego treść jest przedstawiona na listingu 4.34. Aby zaktualizować bazę danych do nowego schematu należy uruchomić polecenie `Model.migrations().migrate()`. OMC spróbuje wykonać migrację, ale nie uda się to. Apache Cassandra nie potrafi zrzutować wartości pól tekstowych na typ `timeuuid`.

```
class User(Model):
    id = TextField(partitioning_key=True)
    # id = UuidField(type=TimeUuid, partitioning_key=True)
    name = TextField()
    # surname = TextField()
```

Rysunek 4.32: Model użytkownika. Linie poprzedzone znakiem # oznaczają zmiany w kolejnej wersji modelu.

```
class MigrationDescription(Migration):
    dependencies = ['__first__']
    operations = [
        Migration.add_model(name='User',
            fields = [
                ('id', TextField(partitioning_key=True)),
                ('name', TextField())
            ])
    ]
```

Rysunek 4.33: Migracja początkowa. Utworzenie modelu danych użytkownika.

```
class MigrationDescription(Migration):
    dependencies = ['initial']
    operations = [
        Migration.alter_field(model_name='User',
            name='id',
            field=UuidField(type=TimeUuid, partitioning_key=True)),
        Migrations.add_field(model_name='User',
            name='surname',
            field=TextField())
    ]
```

Rysunek 4.34: Pierwsza migracja. Do modelu dodano nowe pola oraz zmieniono typ identyfikatora.

Aby migracja zadziałała poprawnie należy nieco zmodyfikować kod z listingu 4.34:

- Zastąpić wpis `alter_field` przez `rename_field`, aby zmienić nazwę kolumny na tymczasową (przykładowo: `temp_id`).
- Dodać wpis `add_field`, który utworzy nowe pole odpowiadające identyfikatorowi typu `timeuuid`.
- Dodać wpis `run_method`, a następnie zaimplementować metodę, która rzutuje wartości z poprzedniej kolumny na nową.
- Dodać wpis `delete_field`, który usunie pole tymczasowe `temp_id`.

4.8.1 Wydajność migracji danych

Migracja danych wymaga wykonania wielu sekwencyjnych zapytań do bazy danych. Nie jest to rozwiązanie efektywne czasowo. W celu zbadania wydajności mechanizmu migracji Autor posłużył się modelem przedstawionym na listingu 4.35.

```
class User(Model):                                class Item(Model):
    id = TextField(key=True)                        id = TextField(key=True)
    name = TextField()                             name = TextField()
    item = DenormalizedField(                      price = TextField()
        relates=Item,                             # color = TextField(
        fields=['name', 'price'])                 default='Black')
```

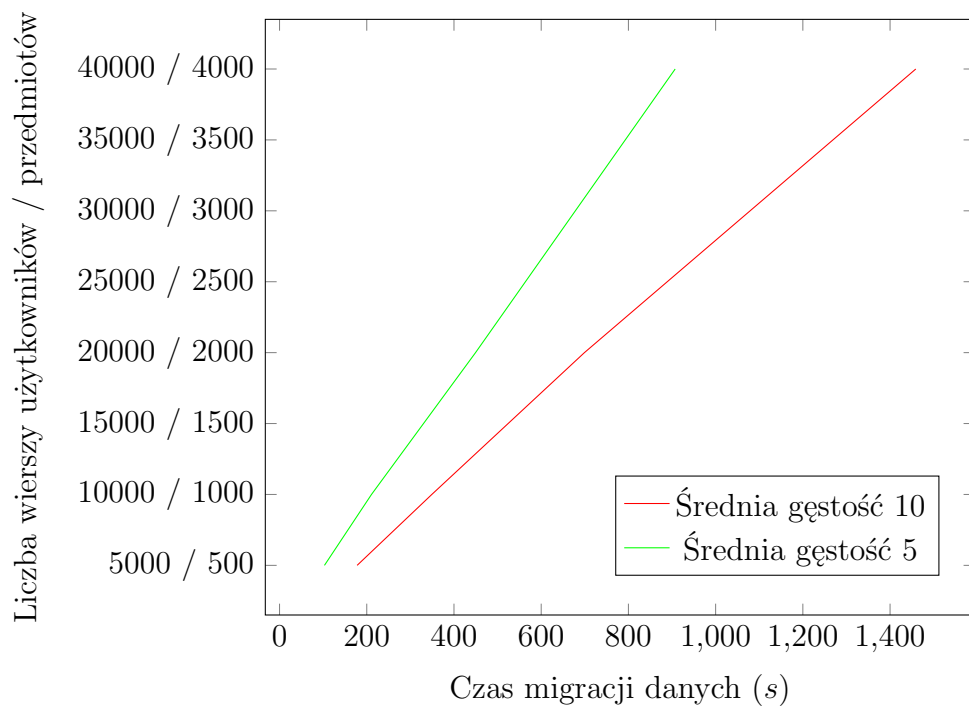
Rysunek 4.35: Definicja modelu danych do testów mechanizmu migracji.

Model składa się z dwóch encji - przedmiotu i użytkownika, który przez zdenormalizowane pola jest związany z przedmiotem. W drugiej wersji modelu do encji `Item` zostało dodane pole `color`. Dodatkowo pole to zostało włączone do zdenormalizowanych pól modelu `User`. W skrypcie migracyjnym ustawiono, że dla istniejących przedmiotów pole `color` ma być uzupełniane losową wartością ze zbioru: czerwony, zielony lub niebieski. Dla tak zdefiniowanego problemu w trakcie migracji muszą zostać wykonane następujące czynności:

1. Do tabeli `item` dodawana jest kolumna `color` typu `text`.

2. Wybierane są wszystkie przedmioty. Dla każdego przedmiotu aktualizowana jest wartość pola `color`.
3. Do tabeli `user` dodawana jest kolumna `item_color` typu `text`.
4. Dla każdego użytkownika pobierane są wszystkie przedmioty. Dla każdego przedmiotu aktualizowana jest wartość zdenormalizowanego pola `item_color`.

Badanie wydajności migracji zostało przeprowadzone na środowisku lokalnym. Zmianie podlegała liczba użytkowników i przedmiotów. Testy przeprowadzone były dla dwóch średnich gęstości zależności: 5 oraz 10. Wyniki badania zostały przedstawione na wykresie 4.36.

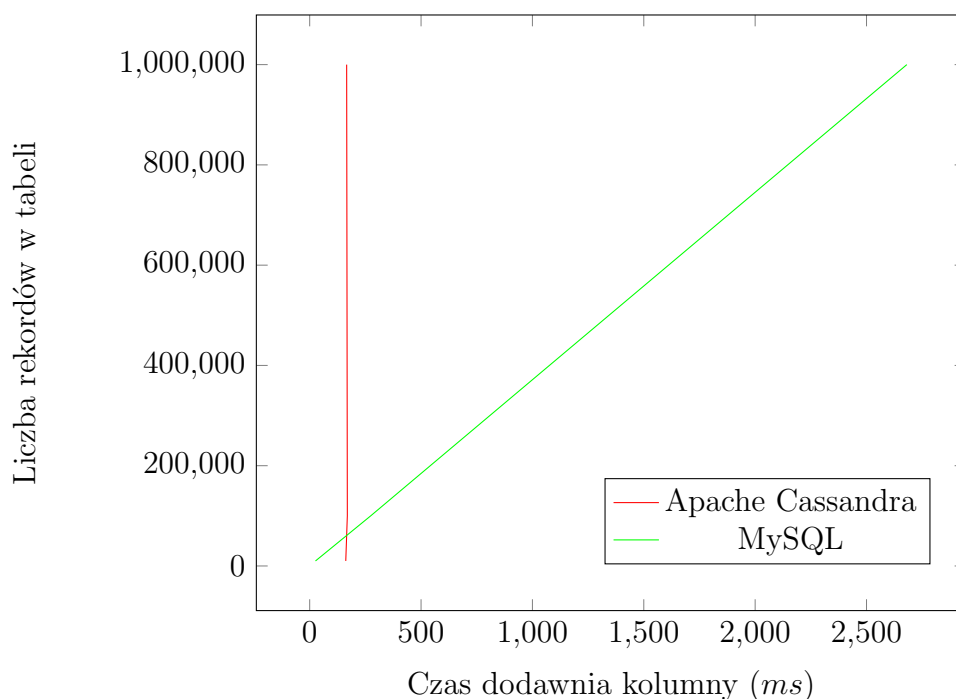


Rysunek 4.36: Porównanie czasów migracji danych dla różnej ilości rekordów.

Z wykresu widać, że migracja danych denormalizacyjnych jest procesem czasochłonnym. Przetworzenie pojedynczej zmiany, dla relatywnie niewielu

wierszy (< 500000), zajmuje 23 minuty. Z tego względu zdecydowanie odradza się wykonywanie tego typu migracji w trakcie produkcyjnego działania systemu.

Patrząc na wyniki przedstawione na wykresie 4.36 należy pamiętać, aby nie utożsamiać ich z wydajnością bazy danych dla operacji przekształcających schemat danych. Ze względu na charakterystykę modelu, Apache Cassandra dużo szybciej niż relacyjne bazy danych dokonuje przekształceń takich jak dodawanie lub usuwanie kolumn. Na wykresie 4.37 przedstawiono porównanie wydajności operacji `ALTER TABLE ... ADD ...` w bazach MySQL oraz Cassandra.



Rysunek 4.37: Porównanie czasów dodawania kolumny do istniejącej tabeli dla bazy relacyjnej oraz Cassandra.

Na podstawie wykresu można zauważyć, że czas dodawania kolumny do istniejącej tabeli w Cassandrze jest niezależny od liczby wierszy. Inaczej jest w przypadku relacyjnej bazy danych, gdzie czas dodawania kolumny rośnie liniowo wraz z ilością rekordów w tabeli. Z tego względu proste migracje

danych mogą być wykonywane podczas produkcyjnego działania systemu.

4.8.2 Potencjalne problemy migracji

Mechanizm migracji dla Cassandra posiada pewne wady w stosunku do analogicznych implementacji dla relacyjnych baz danych. [23] Kluczowy problem jest związany z brakiem możliwości niezawodnego pobrania wszystkich rekordów z danej tabeli. Pomimo, że w najnowszej wersji sterownika dla języka Python został wprowadzony mechanizm paginacji, [24] nadal może wydarzyć się sytuacja przekroczenia czasu żądania podczas pobierania wierszy. OCM używa paginacji o rozmiarze strony równym 10000 rekordów. W zależności od konfiguracji klastra może być to zbyt wysokie ustawienie. Z drugiej strony zmniejszenie rozmiaru paginacji ma zły wpływ na wydajność mechanizmu.

Drugi problem wiąże się z brakiem transakcyjności Apache Cassandra. Jeżeli migracja nie powiedzie się, podczas gdy którakolwiek ze zmian została wprowadzona do modelu danych, ponowne uruchomienie mechanizmu nie zadziała. Z tego względu zaleca się ostrożne wykorzystanie migracji i poprzedzanie każdej takiej operacji wykonaniem kopii zapasowej danych.

4.9 Profilowanie i testowanie modelu danych

OCM posiada zintegrowane narzędzia umożliwiające profilowanie i testowanie wydajności modelu danych. Najbardziej podstawowym narzędziem jest logowanie zapytań. Pozwala ono przechwytywać wszystkie odwołania do bazy i zapisywać je do pliku wraz z czasem wykonania. Dodatkowo użytkownik OMC może zażądać żeby zapisywane były tylko i wyłącznie zapytania, których wykonanie jest dłuższe niż zadany czas. Przykład uruchomienia takiej funkcjonalności został przedstawiony na listingu 4.38.

```
Model.engine.enable_logging(file='log.txt', longer_than=1000)
```

Rysunek 4.38: Wypisywanie do pliku log.txt wszystkich zapytań, których czas wykonania przekracza 1000 milisekund.

Możliwe jest również wypisywanie najwolniejszych zapytań. W trybie śledzenia niewydajnych zapytań OMC zapisuje zadaną liczbę odwołań do bazy o najdłuższym czasie wykonania, posortowanych zaczynając od najwolniejszego. Przykład wykorzystania tej funkcjonalności został przedstawiony na listingu 4.39.

```
Model.engine.trace_slow_queries(amount=10)
# perform ORM database operations
Model.engine.output_slow_queries(file='slow_queries.txt')
```

Rysunek 4.39: Wypisywanie do pliku `slow_queries.txt` 10 najwolniejszych odniesień do bazy danych.

4.9.1 Zasilanie danymi

Nie da się zmierzyć wydajności modelu danych przy pustej bazie, stąd OMC udostępnia dwa narzędzia ułatwiające populację testowymi danymi:

Eksport i import Pozwala na przenoszenie danych pomiędzy instancjami aplikacji. Przykładowym scenariuszem jest przeniesienie danych z serwera produkcyjnego na środowisko testowe w celu optymalizacji modelu.

Generacja Wygenerowanie testowego zestawu danych na podstawie automatycznych lub zdefiniowanych przez użytkownika algorytmów.

Eksport i import wykorzystuje natywny mechanizm języka CQL - polecenia `COPY TO` oraz `COPY FROM`. Służą one do zapisywania danych z Cassandra w formacie CSV¹⁴.

Mechanizm eksportu w OMC wywołuje polecenie `COPY TO` dla wszystkich tabel, a następnie łączy wyniki w jeden plik, w którym definicje tabel znajdują się w wierszach rozpoczynających się znakiem komentarza `#`. Mechanizm

¹⁴Comma Separated Values (ang. „wartości oddzielone przecinkami”) - format pliku, w którym każda linia oznacza wiersz danych, a wartości poszczególnych kolumn są rozdzielone przecinkami.

importu rozdziela pliki, a następnie na każdym z nich wywołuje zapytanie `COPY FROM`. Aby mechanizm działał poprawnie obie instancje Cassandra muszą posiadać już istniejący, spójny model danych. Przykładowe wywołanie zostało przedstawione na listingu 4.40.

```
# first Cassandra instance
Model.engine.export_to_csv(name='data_export')
# second Cassandra instance
Model.engine.import_from_csv(name='data_export')
```

Rysunek 4.40: Przykład eksportowania i importowania danych do/z plików CSV.

Używając mechanizmu eksportu należy pamiętać, że obecna implementacja poleceń `COPY ...` jest bardzo niedopracowana. [22] Problem stanowi wykorzystanie zapytania `SELECT * FROM table WITH LIMIT 99999999`, które służy do pobierania danych. Cechą charakterystyczną dla Cassandra jest ustawiany w konfiguracji węzłów parametr `read_request_timeout_in_ms`, który standardowo ma zbyt restrykcyjną wartość, aby powyższe polecenie mogło wykonać się dla jakichkolwiek danych. Przed uruchomieniem eksportu należy odpowiednio dostosować konfigurację bazy danych.

Niestety, mimo zmiany konfiguracji dla ogromnych zbiorów danych nadal mogą wystąpić przekroczenia czasu żądania. Z tego względu funkcja `export_to_csv` ma zdefiniowany parametr `excludes`, który przyjmuje listę nazw modeli wyłączonych z eksportu. W nowej wersji CQL problemy te mają zostać rozwiązane dzięki wbudowanym w język mechanizmom paginacji.

Mechanizm generacji danych testowych umożliwia definiowanie reguł, według których zostaną stworzone dane. Dzięki temu użytkownik może w łatwy sposób przygotować skrypty zasilające bazę danych w sposób odwzorowujący realne wymagania: długość i objętość danych, gęstość i charakterystykę relacji i tym podobne.

Na listingu 4.41 przedstawiono sposób uruchamiania generacji danych testowych. Parametrem funkcji `populate_test_data` jest słownik, który definiuje nazwy modelu i ilość przypadających mu wierszy. W przypadku po-

minięcia modelu w słowniku dane dla niego nie zostaną wygenerowane.

```
Model.engine.populate_test_data({'User': 10000, 'Address': 1000})
```

Rysunek 4.41: Uruchamianie generacji danych testowych.

Domyślnie używany jest naiwny algorytm generacji danych testowych. Pola zastępowane są wartościami domyślnymi. Przykładowo pola tekstowe uzupełniane są frazą `lorem ipsum`, pola typu `timestamp` lub `timeuuid` wypełniane są czasem obecnym, natomiast zależności generowane są przy założeniu, że istnieje 0.5% szansy na sformułowanie związku.

Dane wygenerowane przy naiwnym podejściu nie nadają się dobrze do testów wydajnościowych. Dlatego użytkownik OMC ma możliwość zdefiniowania własnych algorytmów generacji danych testowych. Wymaga to zdefiniowania klasy wewnętrznej modelu o nazwie `Generator`, a w niej metody `generate()`. Przykład wykorzystania generatora przedstawia listing 4.42.

```
class User(Model):
    username = TextField(partitioning_key=True)
    age = IntegerField()
    class Generator(object):
        def generate(item_no):
            username = 'User: ' + str(item_no)
            age = random.randrange(18, 55)
```

Rysunek 4.42: Przykład wykorzystania generatora danych testowych.

Zarówno generatory, jak również mechanizm importu/eksportu posiadają wiele zastosowań, które nie muszą być ograniczone do badania wydajności. Przykładowe wykorzystania to:

- Migracja danych z innego systemu.
- Synchronizacja danych pomiędzy różnymi środowiskami.
- Zarządzanie kopiami zapasowymi danych.

4.10 Lista życzeń użytkownika - studium przypadku

Encję użytkownika przy pomocy OMC można opisać modelem `User` zaprezentowanym na listingu 4.43. Model posiada trzy pola: identyfikator (`id`), imię (`name`) oraz nazwisko (`surname`).

```
class User(Model):  
    id = UuidField(type=TimeUuid, auto_generate=True, partition_key=True)  
    name = TextField(selectable=True)  
    surname = TextField(selectable=True)
```

Rysunek 4.43: Tabela użytkownika opisana przy pomocy OMC.

Identyfikator jest opisany polem typu `UuidField`. Ze względu na zadeklarowany parametr `type=TimeUuid` pole zostanie zmapowane w Cassandra jako dana o typie `timeuuid`. Mechanizm OMC posiada wbudowaną walidację dla identyfikatorów. Przypisanie polu `id` instancji modelu `User` błędnej wartości spowoduje zgłoszenie wyjątku `AttributeError`. Parametr `auto_generate=True` oznacza, że jeżeli identyfikator nie zostanie ustawiony manualnie, to zostanie wygenerowany automatycznie podczas mapowania modelu na bazę danych. Ustawienie `partition_key=True` oznacza natomiast, że pole będzie częścią klucza wiersza tabeli.

Imię i nazwisko opisane są polem typu `TextField`. Jest ono mapowane na bazodanowy typ `text`. Deklaracja `selectable=True` oznacza, że na polu zostanie utworzony indeks drugiego stopnia. Dzięki temu możliwe jest przeszukiwanie tabeli po wartości tego pola. Na listingu 4.44 przedstawiono odwołanie do OCM, dzięki któremu można wyszukać wszystkich użytkowników o imieniu Jakub i nazwisko Turek.

```
User.objects().find(name='Jakub', surname='Turek')
```

Rysunek 4.44: Przykład wyszukiwania użytkownika po imieniu i nazwisku.

Opis encji przedmiotu w OMC został zaprezentowany na listingu 4.45.


```
class Item(Model):
    id = UuidField(type=TimeUuid, auto_generate=True, partition_key=True)
    name = TextField(searchable=True, length=255)
    price = DecimalField(total_positions=8, decimal_positions=2)
    desc = TextField()
    category = DictionaryField(entries=['BOOKS', 'CLOTHES', 'FURNITURE'])
    weight = DecimalField(total_positions=5, decimal_positions=1)
```

Rysunek 4.45: Tabela przedmiotu opisana przy pomocy OMC.

W stosunku do poprzedniej definicji zostały użyte dwa nowe typy pól: `NumberField` oraz `DictionaryField`. Pierwsze z nich służy do przechowywania wartości liczb dziesiętnych. Cassandra nie posiada odpowiednika znanego z relacyjnych baz danych stałoprzecinkowego typu `NUMBER`. O umieszczanie poprawnie zaokrąglonych wartości musi zadbać autor oprogramowania korzystającego z Apache Cassandra. Aby to ułatwić, OMC dostarcza programowe wsparcie dla operacji stałoprzecinkowych. Pole `DecimalField` posiada dwa parametry. Opcja `total_positions` opisuje maksymalną liczbę cyfr dziesiętnych, z których może składać się liczba, a `decimal_positions` definiuje ile z tych cyfr należy do części ułamkowej pola. `DecimalField` w przezroczysty dla programisty sposób wspiera przypisywanie i wszystkie operacje arytmetyczne dla typów `int`, `long` oraz `float`.

Pole `DictionaryField` dostarcza wsparcia dla słowników jednokrotnego wyboru. Parametr `entries` pozwala na wyspecyfikowanie listy dozwolonych wartości pola. Podanie wartości spoza listy powoduje zgłoszenie wyjątku `AttributeError`. Typ pola jest automatycznie wnioskowany z wartości parametru `entries`. Dla listy elementów typu `str` pole zostanie zmapowane na wartość typu `text`. W przypadku, gdy typ nie może być wywnioskowany na podstawie listy, należy podać go jawnie za pomocą parametru `type`.

Opis encji listy życzeń przedstawia listing 4.46.

Pole `DenormalizedField` spełnia postulaty przedstawione w sekcji 3.5. Istotny jest fakt, że na podstawie parametru `clustering_keys=['itemId']` pole `itemId` staje się częścią złożonego klucza głównego tabeli.

```
class Wishlist(Model):  
    userId = UuidField(type=TimeUuid, partition_key=True)  
    item = DenormalizedField(related=Item,  
                             clustering_keys=['itemId'],  
                             fields=['name', 'price'])
```

Rysunek 4.46: Tabela listy życzeń opisana przy pomocy OMC.

Rozdział 5

Podsumowanie

Tematem pracy było zbadanie możliwości przystosowania istniejących rozwiązań mapowania obiektowo-relacyjnego do bazy danych Apache Cassandra. Celem było zachowanie wysokiej wydajności oraz możliwości stosowania wzorców projektowania, tak aby umożliwić łatwą optymalizację. Już na wstępie pracy okazało się, że taka definicja celów jest wzajemnie sprzeczna. Przeprowadzone badania wykazały, że odwzorowanie relacyjnego modelu danych w Apache Cassandrze jest możliwe, jednakże niweluje większość zalet tego silnika bazodanowego. Uzyskane wyniki wydajności operowania na takim schemacie były porównywalne z relacyjnymi systemami baz danych.

Autor postanowił zrezygnować z jednego spośród wzajemnie sprzecznych założeń i odrzucił pełną zgodność z istniejącymi narzędziami do mapowania obiektowo-relacyjnego. Umożliwiło to znaczne rozszerzenie zakresu pracy. Zamiast prostego odwzorowania dziedziny danych w świecie obiektywnym udało się zaproponować i zrealizować kompleksowy mechanizm wspomagający modelowanie, a także zarządzanie danymi. Rozwiązanie realizuje postawione cele zachowania wydajności i wspomagania wykorzystywania wzorców projektowych. Dodatkowo, przynajmniej częściowo, spełnia postulaty zgodności z mapowaniem obiektowo-relacyjnym dla platformy Django. Integracja z istniejącymi aplikacjami, poza drobnymi modyfikacjami kodu źródłowego, przebiega bezproblemowo.

Praca dowiodła skuteczności podejścia „najpierw kod źródłowy” w odniesieniu do modelowania dziedziny i obsługi danych w Apache Cassandra. Dzięki niemu proces tworzenia i użytkownika wydajnego schematu jest dużo prostszy. Użytkowanie wysokopoziomowych wzorców modelowania nie wymaga zagłębiania się w szczegóły implementacyjne, co ułatwia proces uczenia się. Autor ma nadzieję, że przeprowadzone przez niego badania przyczynią się do rozwoju prac nad modelowaniem dziedziny danych w Apache Cassandrze. W przyszłości rozwiązania wykorzystujące zaprezentowane w pracy pomysły mogą stać się standardem komunikacji dla tej bazy danych.

Praca może być dalej rozwijana w kierunku rozszerzania mechanizmu o kolejne wzorce modelowania. Ponadto implementacja może być dalej optymalizowana pod względem wydajnościowym. Przede wszystkim chodzi o usprawnienie procesów eksportowania i importowania danych z/do Apache Cassandra, a także migracji. Operacje te są bowiem najbardziej czasochłonne spośród wszystkich zaprezentowanych w pracy.

Autor pragnie podziękować Michałowi Aniserowiczowi za inspirację do stworzenia mechanizmów migracyjnych, a także zgodę na wykorzystanie danych osobowych na potrzeby przykładu 4.25. Bez niego nie byłoby możliwe ukończenie badań w pełnym zakresie zaprezentowanym w pracy.

Bibliografia

- [1] *Big Data*. W: Gartner IT Glossary [online], dostęp 25 lipca 2014, <http://www.gartner.com/it-glossary/big-data/>.
- [2] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach Mike Burrows, Tushar Chandra, Andrew Fikes, Robert E. Gruber, *Bigtable: A Distributed Storage System for Structured Data*, Google Inc., 2006.
- [3] Avinash Lakshman, Prashant Malik, *Cassandra - A Decentralized Structured Storage System*, Facebook, 2009.
- [4] *Cassandra Overview*. W: Welcome to Cassandra [online], dostęp 31 lipca 2014, <http://cassandra.apache.org>.
- [5] Dr. Eric A. Brewer, *Consistency vs. Availability (ACID vs. BASE)*, „PODC Keynote”, czerwiec 2000, s. 4-7.
- [6] Patrick McFadin, *The data model is dead, long live the data model*, C* Summit EU, 2013.
- [7] Jay Patel, *Cassandra Data Modeling Best Practices, Part 1*. W: ebay tech blog [online], dostęp 28 lipca 2014, <http://www.ebaytechblog.com/2012/07/16/cassandra-data-modeling-best-practices-part-1/>.
- [8] *API*. W: Cassandra Wiki [online], dostęp 21 lipca 2014, <https://wiki.apache.org/cassandra/API>.

-
- [9] Aleksey Yeschenko, *Cassandra anti-patterns: Queues and queue-like datasets*. W: DataStax Developer Blog [online], dostę 24 lipca 2014, <http://www.datastax.com/dev/blog/cassandra-anti-patterns-queues-and-queue-like-datasets>.
 - [10] Dave Gardner, *Cassandra concepts, patterns and anti-patterns*, Apache-Con EU 2012, listopad 2012, s. 47-49.
 - [11] *User input and Doctrine ORM*. W: Doctrine 2 ORM 2 documentation [online], dostę 30 lipca 2014, <http://doctrine-orm.readthedocs.org/en/latest/reference/security.html>.
 - [12] *Kundera*. W: Github [online], dostę 21 lipca 2014, <https://github.com/impetus-opensource/Kundera>.
 - [13] Patrick McFadin, *Getting Started with Time Series Data Modeling*. W: Planet Cassandra [online], dostę 19 lipca 2014, <http://planetcassandra.org/blog/post/getting-started-with-time-series-data-modeling/>.
 - [14] *About Indexes in Cassandra*. W: Apache Cassandra 1.1 Documentation [online], dostę 23 lipca 2014, <http://www.datastax.com/docs/1.1/ddl/indexes>.
 - [15] Richard Low, *The sweet spot for Cassandra secondary indexing*. W: Richard Low's blog [online], dostę 22 lipca 2014, <http://www.wentnet.com/blog/?p=77>.
 - [16] *BATCH*. W: CQL for Cassandra 1.2 Documentation [online], dostę 23 lipca 2014, http://www.datastax.com/documentation/cql/3.0/cql/cql_reference/batch_r.html.
 - [17] *Using a counter*. W: CQL for Cassandra 2.0 Documentation [online], dostę 21 lipca 2014, http://www.datastax.com/documentation/cql/3.1/cql/cql_using/use_counter_t.html.

- [18] Jeffrey M. Barnes, *Object-Relational Mapping as a Persistence Mechanism for Object-Oriented Applications*, Macalester College, 2007.
- [19] Martin Fowler, *Martin Fowler on ORM Hate*. W: Javalobby [online], dostęp 29 lipca 2014, <http://java.dzone.com/articles/martin-fowler-orm-hate>.
- [20] John O'Conner, *Using the Java Persistence API in Desktop Applications*. W: Oracle Java Platform, Standard Edition [online], dostęp 31 lipca 2014, <http://www.oracle.com/technetwork/articles/javase/persistenceapi-135534.html>.
- [21] *Hibernate OGM roadmap*. W: Hibernate Project [online], dostęp 28 lipca 2014, <http://hibernate.org/ogm/roadmap/>.
- [22] Paul Cannon, *Simple data importing and exporting with Cassandra*. W: DataStax Developer Blog [online], dostęp 24 lipca 2014, <http://www.datastax.com/dev/blog/simple-data-importing-and-exporting-with-cassandra>.
- [23] *Migrations*. W: Django documentation [online], dostęp 27 lipca 2014, <https://docs.djangoproject.com/en/1.7/topics/migrations/>.
- [24] Tyler Hobbs, *DataStax Python Driver 2.0 Released*. W: DataStax Developer Blog [online], dostęp 25 lipca 2014, <http://www.datastax.com/dev/blog/datastax-python-driver-2-0-released>.
- [25] *The write path to compaction*. W: Apache Cassandra 2.0 Documentation [online], dostęp 25 sierpnia 2014, http://www.datastax.com/documentation/cassandra/2.0/cassandra/dml/dml_write_path_c.html.
- [26] Srinath Perera, *Consider the Apache Cassandra database*, IBM developerWorks, czerwiec 2012, s. 4-6.

-
- [27] Billy Bosworth, *Why should I use Cassandra?*. W: DataStax Developer Blog [online], dostęp 14 sierpnia 2014, <http://www.datastax.com/2012/01/why-should-i-use-cassandra>.
- [28] Charles Bell, Mats Kindahl, Lars Thalmann, *MySQL High Availability: Tools for Building Robust Data Centers*, Wyd. 1, Sebastopol, O'Reilly, ISBN 978-0-596-80730-6, s. 165-168.
- [29] *Why Migrate from MySQL to Cassandra?*. W: DataStax White Paper June 2012 [online], dostęp 20 sierpnia 2014, <http://www.datastax.com/wp-content/uploads/2012/08/WP-DataStax-MySQLtoCassandra.pdf>.
- [30] *7 Using Partitioning in an Online Transaction Processing Environment*. W: Oracle® Database VLDB and Partitioning Guide [online], dostęp 27 sierpnia 2014, http://docs.oracle.com/cd/E11882_01/server.112/e25523/part_oltp.htm.
- [31] Paul Gil, *What is 'Saas' (Software as a Service)?*. W: About Technology [online], dostęp 28 sierpnia 2014, http://netforbeginners.about.com/od/s/f/what_is_SaaS_software_as_a_service.htm.
- [32] Mats Kindahl, Alfranio Correia, Narayanan Venkateswaran, *MySQL Sharding: Tools and Best Practices for Horizontal Scaling*, „MySQL Connect”, wrzesień 2013.
- [33] Tyler Hobbs, *Twissandra*. W: GitHub [online], dostęp 29 sierpnia 2014, <https://github.com/twissandra/twissandra>.
- [34] Earl Oliver, *Design and Implementation of a Short Message Service Data Channel for Mobile Systems*, Cheriton School of Computer Science, University of Waterloo, 2007.
- [35] *UUID and timeuuid*. W: CQL for Cassandra 1.2 [online], dostęp 29 sierpnia 2014, http://www.datastax.com/documentation/cql/3.0/cql/cql_reference/uuid_type_r.html.

-
- [36] *Compound keys and clustering*. W: CQL for Cassandra 1.2 [online], dostę-
p 29 sierpnia 2014, http://www.datastax.com/documentation/cql/3.0/cql/ddl/ddl_compound_keys_c.html.
- [37] *CassandraLimitations*. W: Cassandra Wiki [online], dostę-
p 28 sierpnia 2014, <http://wiki.apache.org/cassandra/CassandraLimitations>.
- [38] *Using a composite partition key*. W: CQL for Cassandra 1.2 [online], do-
stęp 30 sierpnia 2014, http://www.datastax.com/documentation/cql/3.0/cql/cql_reference/create_table_r.html.

OŚWIADCZENIE

Oświadczam, że Pracę Dyplomową pod tytułem „Mechanizm modelowania danych i mapowania obiektowego dla Apache Cassandra”, którą kierował dr inż. Jakub Koperwas, wykonałem samodzielnie, co poświadczam własnoręcznym podpisem.

.....

Jakub Turek