

Magazyny danych – stan obecny i kierunki rozwoju

Robert WREMBEL, Zbyszko KRÓLIKOWSKI, Mikołaj MORZY

Politechnika Poznańska, Instytut Informatyki
ul. Piotrowo 3a, 60-965 Poznan

Otrzymano [...]

Streszczenie. Magazyny danych są bardzo dużymi bazami danych, w których gromadzi się dane pochodzące z wielu heterogenicznych źródeł, np. innych scentralizowanych lub rozproszonych baz relacyjnych, relacyjno–obiektowych, obiektowych oraz ze źródeł innych niż bazy danych, takich, jak arkusze kalkulacyjne, pliki XML, zasoby WWW. Potrzeba gromadzenia takich danych jest szczególnie widoczna w sektorze handlowym. Analiza informacji o rynku, popycie i sprzedaży danych towarów, trendach i anomaliach jest jednym z etapów podejmowania strategicznych decyzji w przedsiębiorstwie. Analizę przeprowadza się zwykle za pomocą aplikacji OLAP wykorzystujących bardzo złożone zapytania i operujących na ogromnych ilościach danych, rzędu setek, czy tysięcy gigabajtów. Kluczowym problemem w magazynach danych jest więc efektywność przetwarzania informacji. W ostatnim czasie wiele uwagi poświęcono opracowaniu modeli danych (wielowymiarowy model danych), metod fizycznej implementacji modelu wielowymiarowego (relacyjny i wielowymiarowy OLAP), nowych algorytmów optymalizacji zapytań (indeksy bitmapowe).

Artykuł ten jest poświęcony omówieniu procesu konstruowania magazynu danych, zagadnień związanych z pielęgnacją i odświeżaniem magazynu oraz kwestiom architektury magazynu danych. W artykule przedstawiono nowy model przetwarzania (OLAP), wielowymiarowy model danych, metody jego implementacji oraz opisano podstawy zagadnienia optymalizacji zapytań w środowisku magazynu danych.

Słowa kluczowe : magazyn danych, OLAP, OLTP, analiza danych, wielowymiarowy model danych, kostka danych, indeks bitmapowy

1 Wprowadzenie

Sposób, w jaki użytkownik korzysta z bazy danych nazywamy modelem przetwarzania. Tradycyjny model przetwarzania – *przetwarzanie transakcji w trybie on-line* (ang. On-Line Transaction Processing OLTP) doskonale nadaje się do obsługi bieżącej działalności firmy, obsługi operacji bankowych i księgowych, wykonywania powtarzalnych i dobrze zdefiniowanych procesów. Model ten charakteryzuje się krótkimi i prostymi transakcjami które operują na niewielkich ilościach danych. Większość operacji stanowią operacje modyfikacji i zapisu krotek, równoległe wykonują się setki bądź tysiące transakcji, stąd kluczowe znaczenie mają kwestie integralności, odtwarzania czy izolacji. Użytkownikami systemów OLTP są najczęściej urzędnicy, pracownicy obsługi klienta, itp. Niestety, model OLTP nie znajduje zastosowania w procesach analizy danych, w systemach wspomagania decyzji czy też w przypadku procesów nie do końca dobrze zdefiniowanych.

Wspomaganie decyzji stanowi drugi obszar zastosowania systemów baz danych. Rozwijające się galezie przemysłu, handlu, medycyny, nauki wymagają składowania i przetwarzania ogromnych ilości danych. Dane są zwykle przechowywane w systemach informatycznych posiadających różne struktury i wykorzystujących różne modele danych (np. hierarchiczne, relacyjne, obiektowe), w dokumentach tekstowych, czy arkuszach kalkulacyjnych. Ponieważ przedsiębiorstwa, instytucje, organizacje (producenci danych) są często geograficznie rozproszone, więc same dane mają również charakter rozproszony. Posiadanie danych opisujących działanie przedsiębiorstwa w dłuższym przedziale czasu pozwala na analizę trendów, anomalii, poszukiwania wzorców zachowań (wzorców zakupów, podobieństw między klientami, itp.) czy też odkrywania wiedzy. Dane przechowywane w bazie danych zawierają olbrzymią ilość potencjalnie użytecznej wiedzy, która może zostać użyta w procesie podejmowania decyzji strategicznych dotyczących działalności przedsiębiorstwa. Dla potrzeb analizy danych opracowano zatem nowy model przetwarzania, nazywany *przetwarzaniem analitycznym w trybie on-line* (ang. On-Line Analytical Processing OLAP) oraz stworzono nowy typ architektury bazy danych nazywany *magazynem danych* (ang. data warehouse). Magazyny danych są bardzo dużymi bazami danych, w których gromadzi się dane pochodzące z wielu heterogenicznych źródeł, np. innych scentralizowanych lub rozproszonych baz relacyjnych, relacyjno–obektowych, obiektowych oraz ze źródeł innych niż bazy danych, np. arkusze kalkulacyjne, dokumenty XML lub pliki tekstowe. Analiza informacji o rynku, popycie i podażu danych towarów, trendach, modach i anomaliami jest bardzo istotna dla procesu podejmowania strategicznych decyzji w przedsiębiorstwie. Analizę przeprowadza się zwykle poprzez bardzo złożone zapytania operujące na ogromnych ilościach danych, rzędu setek, czy tysięcy gigabajtów. Kluczowym problemem w magazynach danych jest więc efektywność przetwarzania informacji.

W niniejszym artykule przedstawiony zostanie szczegółowo model przetwarzania OLAP oraz metody modelowania, projektowania i implementacji magazynu danych.

2 Model przetwarzania analitycznego w trybie on-line (OLAP)

Komercyjnie dostępne systemy transakcyjne (systemy zarządzania bazami danych SZBD) dostarczają efektywnych rozwiązań dla takich problemów jak: efektywne i bezpieczne przechowywanie danych, transakcyjne odtwarzanie danych, dostępność danych, optymalizacja dostępu do danych, zarządzanie współbieżnością. W znacznie mniejszym stopniu systemy te wspomagają operacje agregacji danych, wykonywania podsumowań czy też optymalizacji złożonych zapytań formułowanych ad hoc. Systemy te w niewielkim stopniu wspomagają również integracje danych z różnych heterogenicznych źródeł danych. W ostatnim czasie prace badawcze i rozwojowe prowadzone nad rozszerzeniem funkcjonalności systemów baz danych doprowadziły do opracowania nowego modelu przetwarzania danych, którego podstawowym celem jest wspomaganie procesów podejmowania decyzji, oraz opracowania nowego typu relacyjnej bazy danych nazwanego *magazynem danych* (ang. data warehouse).

Nowy model przetwarzania danych, nazwany *przetwarzaniem analitycznym on-line* (ang. On-Line Analytical Processing OLAP), ma za zadanie wspieranie procesów analizy magazynów danych dostarczając narzędzi umożliwiających analizie magazynu w wielu „wymiarach” definiowanych przez użytkowników (czas, miejsce, klasyfikacja produktów, itp.). Analiza magazynu polega na obliczaniu agregatów dla zadanych „wymiarów” magazynu. Należy podkreślić, że proces analizy jest całkowicie sterowany przez użytkownika. Mówimy czasami o *analizie danych sterowanej zapytaniami* (ang. query-driven exploration). Typowym przykładem takiej analizy jest zapytanie o sprzedaż produktów w supermarkecie w kolejnych kwartałach, miesiącach, tygodniach, itp., zapytanie o sprzedaż produktów z podziałem na rodzaje produktów (AGD, produkty spożywcze, kosmetyki, itp.), czy wreszcie zapytanie o sprzedaż produktów z podziałem na oddziały supermarketu. Odpowiedzi na powyższe zapytania umożliwiają decydentom określenie wąskich gardeł sprzedaży, produktów przynoszących deficyt, oraz podjęcie odpowiednich działań poprawiających sytuację. Niestety, obliczenia agregatów są bardzo czasochłonne. Częściowym rozwiązaniem problemu efektywności analizy magazynu danych jest wstępne przeprowadzenie obliczeń i zmaterializowanie otrzymanych wyników w magazynie danych w celu ich późniejszego wykorzystania. Jednakże, liczba możliwych agregatów jest wykładniczo zależna od liczby „wymiarów”, stąd rodzi się pytanie, które z agregatów zmaterializować, które agregaty pozostawić do obliczeń w trybie on-line, oraz w jaki sposób wykorzystywać zmaterializowane agregaty do obliczenia innych agregatów. Pojawia się szereg dodatkowych bardzo interesujących pytań: w jaki sposób reprezentować i realizować dostęp do agregatów, w jaki sposób aktualizować agregaty, czy nie należy również zmaterializować pośrednich wyników obliczeń (nie tylko samych agregatów), np. wyników niektórych operacji połączeń, które są wspólne dla wielu agregatów.

Model OLAP adresowany jest przede wszystkim do decydentów, menadżerów, analityków systemowych. Model ten charakteryzuje się m.in. małą liczbą złożonych i kosztownych zapytań.

Znakomita większość transakcji w modelu OLAP składa się z operacji odczytu dużych wolumenów danych. Dane są historyczne, zagregowane, zintegrowane i wielowymiarowe. Pojedyncze zapytanie operuje na tysiącach bądź milionach rekordów. Podstawowa miara wydajności systemów nie jest ilości transakcji na jednostkę czasu (jak w przypadku systemów OLTP), lecz czas odpowiedzi na pojedyncze zapytanie. Ponieważ operacja modyfikacji magazynu danych jest bardzo rzadka, kryteria izolacji i integralności transakcji nie są tak istotne, jak w przypadku systemów OLTP.

Analiza danych ma olbrzymie znaczenie dla procesów wspomaganie podejmowania decyzji. Jednakże, aby przeprowadzić taką analizę, należy dysponować odpowiednimi danymi opisującymi działalność przedsiębiorstwa. Bardzo rzadko informacje te są dostępne w jednej bazie danych. Z reguły są one rozproszone po wielu oddziałowych, rozproszonych geograficznie i heterogenicznych bazach danych. Stąd, opracowując koncepcje systemu wspomaganie podejmowania decyzji należy odpowiedzieć na dwa zasadnicze pytania odnośnie architektury takiego systemu i modelu przetwarzania. Pierwsze pytanie brzmi: czy analiza powinna mieć charakter rozproszony czy scentralizowany, innymi słowami, czy dane należy zgromadzić i przetwarzać w jednym miejscu w sposób scentralizowany, czy też korzystając z mechanizmu transakcji rozproszonych można przetwarzać dane w sposób rozproszony. Drugie pytanie dotyczy koegzystencji dwóch systemów – systemu bieżącej obsługi działania przedsiębiorstwa oraz systemu wspomaganie podejmowania decyzji. Oba systemy operują na tych samych danych, stąd pytanie, czy oba modele OLAP i OLTP mogą współistnieć w tym samym systemie bazy danych, czy też powinny funkcjonować niezależnie.

W chwili obecnej odpowiedź na powyższe pytania brzmi: analiza powinna mieć charakter scentralizowany, a modele OLAP i OLTP powinny funkcjonować niezależnie. Oczywiście, odpowiedź na pytania o architekturę i model przetwarzania jest uzależniona od aktualnego stanu rozwoju technologii informatycznej. Ze względu na charakter i prędkość obliczeń, częściowo również ze względu na problem autoryzacji dostępu do danych, analiza danych jest aktualnie prowadzona w sposób scentralizowany. Wraz z rozwojem sieci komputerowych, wzrostem prędkości transmisji danych, należy się jednak spodziewać przechodzenia od modelu przetwarzania analitycznego scentralizowanego do modelu przetwarzania analitycznego rozproszonego.

3 Projektowanie magazynu danych

Magazyn danych jest „(...) zorientowana tematycznie, zintegrowana, zmienna w czasie i nieulotna kolekcja danych wspierających proces wspierania decyzji” (W.H.Inmon). Jest to zbiór technologii, których celem jest wspieranie decydentów (menedżerów, analityków) i umożliwianie im podejmowania lepszych i szybszych decyzji. Zakłada się, że w magazynie danych znajdują się właściwe dane we właściwym miejscu, właściwym czasie i właściwym koszcie dla właściwych decyzji.

Magazyny danych są **zorientowane tematycznie** w zależności od podstawowego obszaru działalności danego przedsiębiorstwa. Dla przykładu, operacyjne bazy danych firmy ubezpieczeniowej mogą być zorientowane na poszczególne segmenty działalności firmy, np. ubezpieczenia na życie, ubezpieczenia samochodów, domów, itp. Magazyn danych tej firmy będzie zorientowany na klientów, typy ubezpieczeń, polisy, konta, zadania wypłat, itp. Należy raz jeszcze powtórzyć, że celem magazynu danych nie jest wspieranie bieżącej działalności danego przedsiębiorstwa, lecz asystowanie decydentom w analizie rynku, w podejmowaniu decyzji strategicznych dla przedsiębiorstwa, w znajdowaniu wąskich gardeł w działalności firmy.

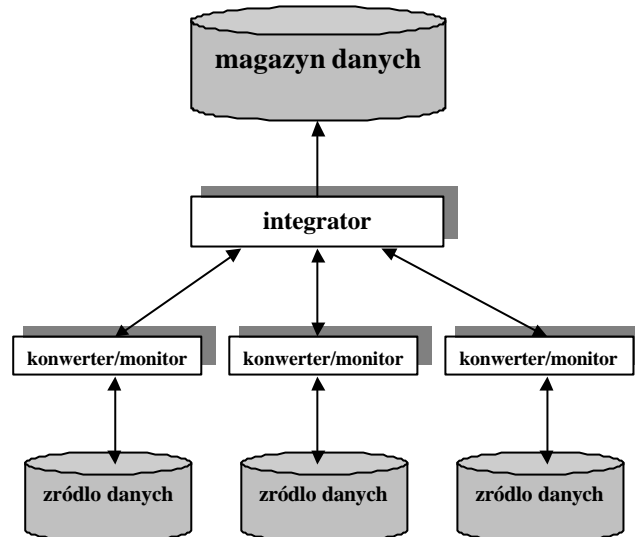
Aby wypełnić postawione przed nimi cele magazyny danych muszą dysponować pełną wiedzą dotyczącą całości funkcjonowania firmy. Niezbędna staje się więc **integracja danych** z wielu heterogenicznych źródeł. Ważnym wymaganiem systemów informatycznych jest dostęp do informacji pochodzących zarówno z heterogenicznych baz danych (wykonanych w różnych technologiach), jak i ze źródeł innych niż bazy danych, np. arkusze kalkulacyjne, dokumenty tekstowe, pliki semi- i strukturalne HTML, SGML czy XML. Informacje przechowywane w różnych źródłach są integrowane w magazynie danych. Integracja polega na rozwiązywaniu niekompatybilności między różnymi typami danych, uspoźnianiu nazewnictwa, konwersji wartości, itp.

Dane operacyjne są regularnie aktualizowane i zmieniane. Transakcje w systemach OLTP dokonują dostępu do pojedynczych rekordów, zmieniają ich wartości, usuwają bądź wstawiają rekordy. Magazyny danych natomiast są **nieulotne**, po załadowaniu danych do magazynu dane te podlegają tylko operacji odczytu. Aktualizacja danych przechowywanych w magazynie danych, czyli *odświeżenie danych* (ang. refresh) odbywa się w trybie wsadowym w określonych odstępach czasu (raz na tydzień, raz na miesiąc).

Istotną cechą magazynów danych jest ich **zmiennosc w czasie**. Horyzont czasowy magazynu danych jest znacząco większy niż horyzont czasowy operacyjnych baz danych. Bazy operacyjne przechowują aktualne wartości danych i nie zawsze zawierają element czasu. W przeciwieństwie do nich, magazyny danych przechowują całą historię danych (czyli zbiór migawek zrobionych w pewnych odstępach czasowych) i czas stanowi zawsze jeden z podstawowych elementów składowych magazynu danych.

Technologia magazynów danych stała się istotną strategią integracji heterogenicznych źródeł informacji. Zgodnie z raportem grupy META wydatki na budowę magazynów danych osiągnęły w roku 1998 wysokość \$15 mld i z każdym rokiem wzrastają. Szacunkowy koszt jednego projektu związanego z tworzeniem bądź modernizacją magazynu danych wynosi obecnie \$3 mln i również wzrasta.

Architektura magazynu danych została przedstawiona na Rys. 1 [16]. Obiekty oznaczone jako *źródło danych* reprezentują wymienione wyżej heterogeniczne źródła informacji. Z każdym z takich źródeł jest związana warstwa oprogramowania –*konwerter/monitor*.



Rys.1. Architektura magazynu danych

Zadaniem modułu *konwertera* jest transformowanie danych z formatu wykorzystywanego w źródle, do formatu wykorzystywanego w magazynie. Dlatego, dla każdego modelu danych źródłowych konieczne jest zastosowanie specyficznego modułu *konwertera*. Przykładowo, jeśli źródło przechowuje dane w dokumentach tekstowych, a magazyn został zaprojektowany z wykorzystaniem modelu relacyjnego, to *konwerter* musi zapewnić poprawne odwzorowanie danych z plików w strukturę modelu relacyjnego.

Konwersja danych ze źródeł do magazynu rozpoczyna się od procesu *ekstrakcji danych* (ang. data extraction). Dane są czytane ze struktur źródłowych (plików tekstowych, arkuszy kalkulacyjnych, sieci WWW, baz danych) za pomocą *bramek* (ang. gateways) i standardowych interfejsów (Information Builders EDA/SQL, ODBC, JDBC, Oracle Open Connect, Sybase Enterprise Connect, Informix Enterprise Gateway, itp.), choć w wielu przypadkach występuje potrzeba zaimplementowania wyspecjalizowanych procedur do ekstrakcji danych ze źródeł niestandardowych. Kolejnym krokiem jest *czyszczenie danych* (ang. data cleaning, data cleansing, data scrubbing). Proces ten ma na celu zapewnienie jakości i poprawności danych. Ma on szczególne znaczenie w przypadku, gdy źródła danych są wysoce heterogeniczne, posługują się różnymi schematami danych lub opisują te same obiekty świata rzeczywistego w różny sposób. Dane przed transformacją często nazywane są danymi „brudnymi”. Przykładami zabrudzenia danych są [5]:

- różne formaty danych dla tego samego pola (np. informacja o województwie może być przedstawiona jako pełna nazwa, skrót nazwy lub specjalny identyfikator)
- niespójne wartości tej samej danej spowodowane błędami przy wprowadzaniu
- niezgodności między wartością atrybutu i jego nazwą (np. pole *Nazwa* może zawierać nazwę firmy lub nazwisko indywidualnego klienta)
- brakujące wartości, które zgodnie ze schematem magazynu danych powinny być wypełnione
- redundantne informacje na temat jakiegoś obiektu świata rzeczywistego (mogą pojawiać się zarówno w ramach jednego źródła informacji, jak i w dwóch różnych źródłach)

Podstawowymi metodami czyszczenia danych są:

- konwersja i normalizacja: transformacja i standaryzacja heterogenicznych formatów danych
- czyszczenie specjalne: uspojnianie wartości pola na podstawie słownika synonimów
- czyszczenie niezależne od domeny: porównywanie pól pochodzących z różnych źródeł w celu określenia stopnia ich podobieństwa
- czyszczenie oparte na regułach: określanie podobieństwa między polami na podstawie zespołu predefiniowanych reguł

Po wyczyszczeniu danych następuje etap *ładowania danych* (ang. data loading) do magazynu danych. Procesem ładowania zarządza zaznaczony na Rys.1. modul *integratora*. Ładowanie danych pociąga za sobą dodatkowe przetwarzanie, np. sprawdzanie ograniczeń integralnościowych, sortowanie, podsumowywanie, budowanie indeksów, itp. Ładowanie danych odbywa się najczęściej w trybie wsadowym. Aplikacja dokonująca załadowania danych musi pozwalać administratorowi na monitorowanie procesu ładowania, zawieszanie i odwieszanie ładowania, restartowanie po awarii (np. po naruszeniu ograniczeń integralnościowych). Proces ładowania danych może być bardzo czasochłonny i trwać wiele godzin bądź dni. W rzeczywistości taki proces może być traktowany jako jedna transakcja konstruująca nową bazę danych. Po załadowaniu danych wszystkie następujące zmiany w danych źródłowych są propagowane do magazynu danych podczas odświeżania. *Odświeżanie danych* (ang. data refresh) to proces propagowania zmian zachodzących w źródłach danych do magazynu. Z odświeżaniem wiąże się dwa istotne problemy: w jaki sposób odświeżać oraz kiedy odświeżać. Odświeżanie magazynu danych może zachodzić okresowo (np. pod koniec dnia, kiedy magazyn danych nie jest używany), natychmiastowo (zmiany propagowane są do magazynu natychmiast po wystąpieniu w danych źródłowych) lub w sposób określony przez charakter źródła (np. po wystąpieniu określonej ilości zmian). Jeśli chodzi o metody odświeżania można wyróżnić zasadniczo dwie metody:

- data shipping: relacja w magazynie danych jest traktowana jak zdalna migawka źródła danych. Wyzwalacze działające w źródle są odpowiedzialne za aktualizacje logu migawki i propagację danych do magazynu
- transaction shipping: log transakcji w źródle danych jest okresowo przeglądany celem wykrycia zaistniałych zmian. Wszelkie zmiany są przysyłane do serwera replikacji który przysyła nowe transakcje do magazynu danych

Zadaniem modułu *monitora* jest wykrywanie zmian w danych źródłowych i ich przekazywanie do warstwy oprogramowania *integratora* (po uprzedniej konwersji do modelu danych magazynu). Sposób wykrywania zmian w danych źródłowych zależy od własności samych źródeł. Wyróżnia się cztery następujące rodzaje źródeł danych [16]:

- aktywne (ang. active sources), tzn. posiadające zaimplementowane mechanizmy wyzwalaczy, które informują *monitor* o zmianach zachodzących w danych źródłowych
- utrzymujące dzienniki operacji wykonywanych na danych źródłowych (ang. logged sources) – zmiany są wykrywane przez analizę zawartości dziennika przez moduł *monitora*
- umożliwiające wydawanie zapytań (ang. queryable sources) – w celu wykrycia zmian w danych źródłowych *monitor* okresowo wydaje zapytania do źródła
- wspierające mechanizm migawek (ang. snapshot sources) – w tego rodzaju źródłach zmiany wykrywa się przez porównanie zawartości kolejnych migawek

W przypadku gdy problem wykrywania i propagowania zmian w danych źródłowych jest trudny do rozwiązania, stosuje się okresowe uaktualnianie całego magazynu w trybie wsadowym. Wówczas magazyn jest niedostępny dla użytkowników. Ta technika jest wykorzystywana również w przypadku gdy aktualność danych nie jest konieczna dla poprawnego działania firmy, instytucji czy organizacji, oraz gdy dopuszczalna jest czasowa niedostępność danych.

Istotnym składnikiem magazynu danych jest *repozytorium metadanych* (ang. metadata repository), w którym przechowywane są informacje wspomagające zarządzanie magazynem. Metadane wspomagają też efektywny dostęp do zasobów magazynu poprzez przechowywanie informacji o zawartości magazynu, zależnościach między poszczególnymi komponentami systemu lub ich fizycznej lokalizacji w obrębie magazynu. Repozytorium zazwyczaj zawiera następujące informacje:

- listę źródłowych baz danych i opis ich zawartości
- opisy i charakterystyki bramek między bazami źródłowymi a magazynem
- schemat magazynu
- definicje perspektyw i danych wyliczalnych, przechowywanych w magazynie

- opisy wymiarów i hierarchii (patrz dalej)
- predefiniowane zapytania i raporty
- lokalizacje tematycznych hurtowni danych
- zasady czyszczenia, transformacji i korekcji danych źródłowych
- zasady odswieżania danych
- profile użytkowników i profile grup użytkowników
- dane dotyczące bezpieczeństwa magazynu (autoryzacja użytkowników, prawa dostępu do poszczególnych komponentów magazynu)

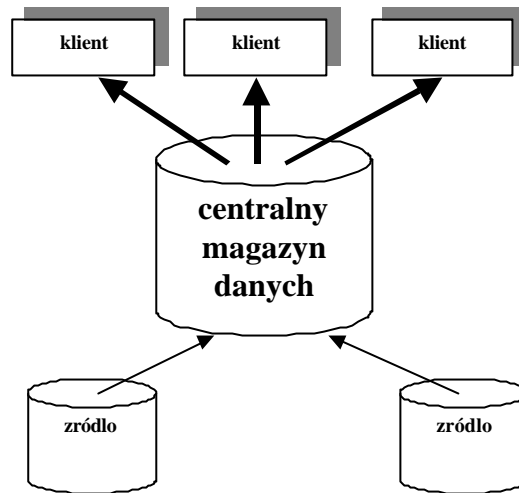
Istniejące komercyjne systemy zarządzania bazami danych wspierające technologie hurtowni danych umożliwiają integracje danych pochodzących z systemów relacyjnych, dostarczanych przez największych wytwórców. Firma Oracle dostarcza oprogramowanie o nazwie *Oracle Loader* i *Oracle Transparen Gateways* umożliwiające dostęp oraz eksportowanie i importowanie danych z większości systemów relacyjnych. IBM umożliwia wczytywanie i transformowanie danych do formatu bazy DB2 za pomocą oprogramowania o nazwie *DataPropagator*, *DataRefresher* i *DataHub*. *SourcePoint* jest oprogramowaniem firmy Software AG służącym do ładowania hurtowni danymi pochodzącymi z innych relacyjnych źródeł. Informix realizuje dostęp do danych zgromadzonych w bazie *Oracle*, *Sybase* i *DB2* za pomocą oprogramowania *Informix–Enterprise Gateway*. Oprogramowanie firmy Sybase umożliwia bezpośrednie wczytywanie danych z wielu systemów relacyjnych. Wszystkie z istniejących na rynku systemów wspierają standard ODBC (ang. Open Database Connectivity) dostępu do bazy danych.

4 Fizyczne struktury magazynów danych

Istnieją trzy podstawowe architektury konstrukcji magazynów danych [9]:

- architektura scentralizowana
- architektura sfederowana
- architektura wielowarstwowa

W architekturze scentralizowanej istnieje tylko jeden magazyn danych przechowujący wszystkie dane potrzebne do analiz. Architektura scentralizowana została przedstawiona na Rys.2.

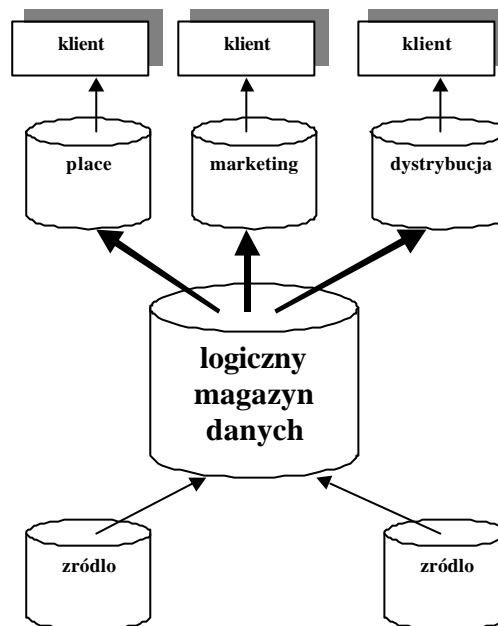


Rys.2. Architektura scentralizowana

Efektywność systemu scentralizowanego jest mniejsza niż systemu rozproszonego. Wszystkie zapytania i operacje odświeżenia magazynu danych wykonywane są w jednej bazie danych. Zaletą architektury scentralizowanej jest łatwość dostępu do danych i łatwość administrowania magazynem. W architekturze scentralizowanej nie zachodzi problem integracji schematów magazynu ponieważ występuje tylko jeden spójny model danych.

Podstawowymi powodami, dla których stosuje się rozproszone architektury magazynów danych, są: skalowalność, zwiększenie dostępności i odporności na awarie oraz zmniejszenie czasu odpowiedzi systemu (dane są zlokalizowane bliżej aplikacji klientów, w tematycznych magazynach danych wolumeny przeszukiwanych danych są mniejsze).

W przypadku decentralizacji magazynu danych każdy z węzłów rozproszonego magazynu danych posiada własne repozytorium metadanych. W architekturze sfederowanej dane są zintegrowane logicznie, lecz fizycznie są przechowywane w osobnych bazach danych. Ważną cechą architektury sfederowanej jest więc fakt, że magazyn danych jest czysto wirtualny. *Tematyczne magazyny danych* (ang. data marts) przechowują szczegółowe dane, które są wykorzystywane do wykonywania zapytań. Schemat architektury sfederalizowanej przedstawiono na Rys.3.



Rys.3. Architektura sfederowana

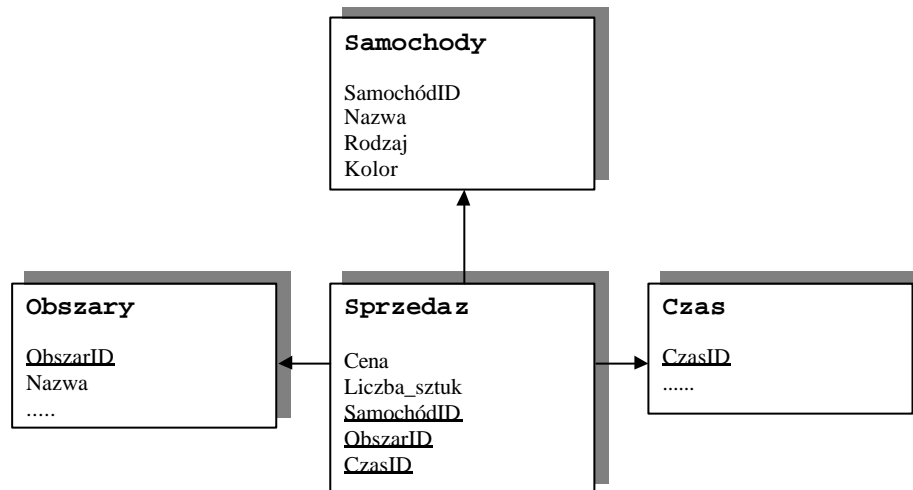
W przypadku architektury wielowarstwowej centralny magazyn danych jest fizyczny. Dodatkowo w lokalnych bazach danych przechowywane są dane charakterystyczne dla danego oddziału. Lokalne bazy danych przechowują kopie tych fragmentów magazynu, które są istotne dla danego oddziału. Kopie te są wstępnie przetworzone i wzbogacone o agregaty, statystyki, itp. Zawartość lokalnych baz danych nie jest jednak tak szczegółowa, jak w przypadku architektury sfederowanej.

5 Modele danych dla magazynów danych

Bazy danych wspierające technologie magazynów danych można podzielić na dwa rodzaje ze względu na wykorzystywane przez nie modele danych. Pierwszy, to magazyny danych wykorzystujące *model relacyjny*, nazywane również *ROLAP* (ang. Relational OLAP), drugi to magazyny danych wykorzystujące wielowymiarowy model danych, nazywane również *MOLAP* (ang. Multidimensional OLAP).

5.1 ROLAP

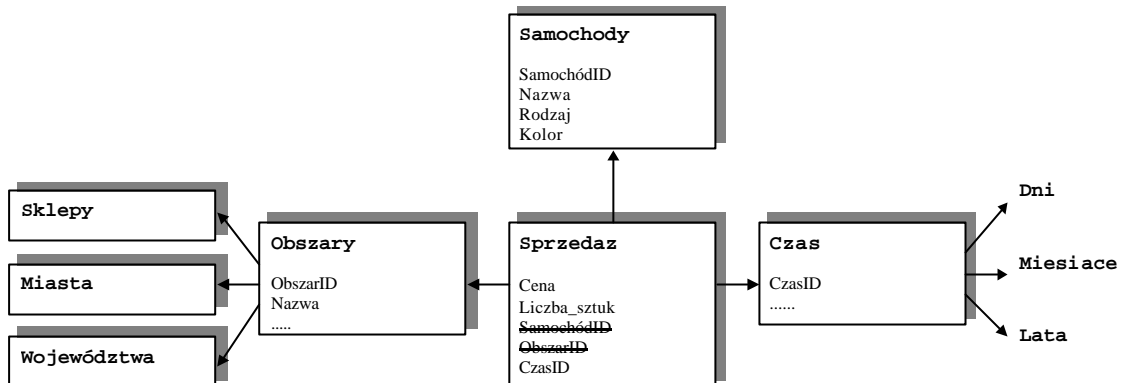
Magazyn danych tego typu zbudowany jest w oparciu o system zarządzania relacyjną bazą danych posiadający mechanizmy efektywnego przetwarzania zapytań typu OLAP. Zwykle schemat takiego magazynu posiada strukturę *gwiazdy* (ang. star schema) lub strukturę bardziej złożoną, przypominającą *platek śniegu* (ang. snowflake schema). W celu skrócenia czasu potrzebnego na wyznaczenie wyników zapytania relacje bazy danych są często denormalizowane, np. zawierają wartości zregulowane, są wynikiem połączenia wielu innych relacji [2]. Przykładowy schemat o strukturze gwiazdy przedstawia Rys.4.



Rys.4. Schemat gwiazdy

Centralna relacja *Sprzedaz* zawiera informacje o sprzedaży samochodów, w pewnych obszarach geograficznych, w określonym czasie. Relacje *Samochody*, *Obszary* i *Czas* są nazywane *wymiarami* (ang. dimensions), natomiast relacja centralna jest nazywana *relacją faktów* (ang. fact table). Atrybuty relacji faktów przechowujące informacje o sprzedaży są nazywane *miarami* (ang. measures), np. *Cena*, *Liczba_sztuk*. Relacja faktów *Sprzedaz* zawiera również atrybuty *SamochódID*, *ObszarID*, *CzasID*, których wartości wskazują na odpowiednie wymiary i pozwalają zidentyfikować dany fakt (sprzedaz pewnej ilości samochodów za pewną cenę) w przestrzeni wymiarów (jaki samochód? gdzie sprzedany? kiedy?).

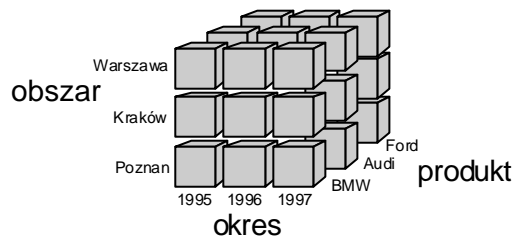
Wymiary mogą tworzyć hierarchie, wówczas schemat przybiera postać płatka śniegu [8]. Rysunek 5 przedstawia przykładowy schemat o takiej strukturze.



Rys.5. Schemat płatka śniegu

5.2 MOLAP

Magazyn danych zaprojektowany w technologii MOLAP do przechowywania danych wykorzystuje wielowymiarowe tablice (ang. multidimensional arrays) zwane też kostkami danych (ang. data cubes). Tablice te zawierają wstępnie przetworzone (m.in. zagregowane) dane pochodzące z wielu źródeł. Przykładowa tablica zawierająca trzy wymiary: *Obszar*, *Okres* i *Produkt* oraz zagregowane informacje o sprzedaży samochodów w poszczególnych latach w wybranych miastach została przedstawiona na Rys.6.



Rys.6. Tablica trójwymiarowa

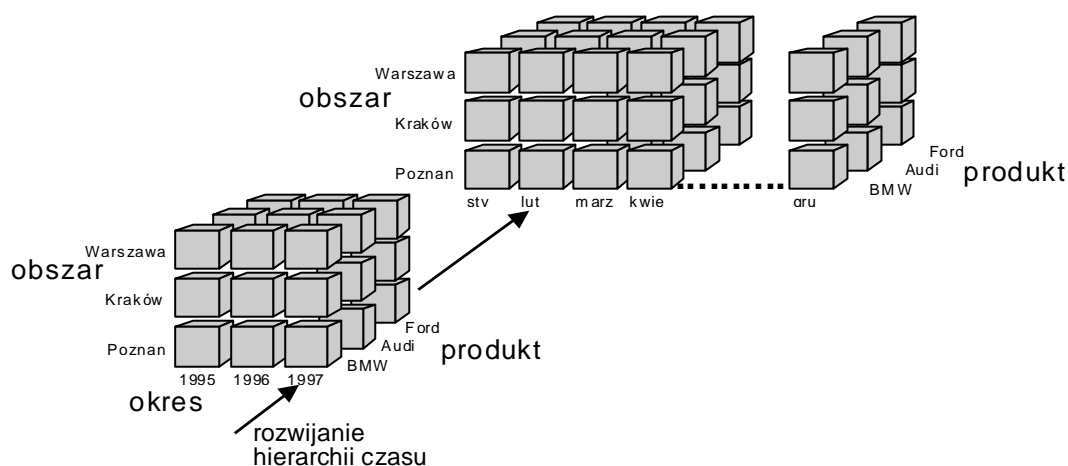
Analize danych wielowymiarowych wspomagają specjalne operatory, do których należą:

- *wyznaczanie punktu centralnego* (ang. pivoting)

Operacja ta polega na wskazaniu miary i określeniu wymiarów, w których wybrana miara będzie prezentowana. Przykładowo, w wymiarze produktu reprezentującego samochód marki „BMW” i wymiarze obszaru reprezentującego sklepy województwa poznańskiego może być prezentowana liczba sprzedanych samochodów.

- *rozwijanie* (ang. drilling down)

Rozwijanie polega na zagłębianiu się w hierarchie danego wymiaru w celu przeprowadzenia bardziej szczegółowej analizy danych. Jako przykład rozważmy informacje o sprzedaży samochodów marek BMW, Audi i Ford, w latach 1995, 1996 i 1997, w poszczególnych miastach (por. Rys. 6). W celu dokonania analizy sprzedaży w poszczególnych miesiącach roku 1997 należy rozwinąć hierarchie reprezentującą czas, tj. rok 1997 (zob. Rys. 7). Analiza sprzedaży w poszczególnych dniach wybranego miesiąca będzie możliwa po rozwinięciu hierarchii reprezentującej ten miesiąc.



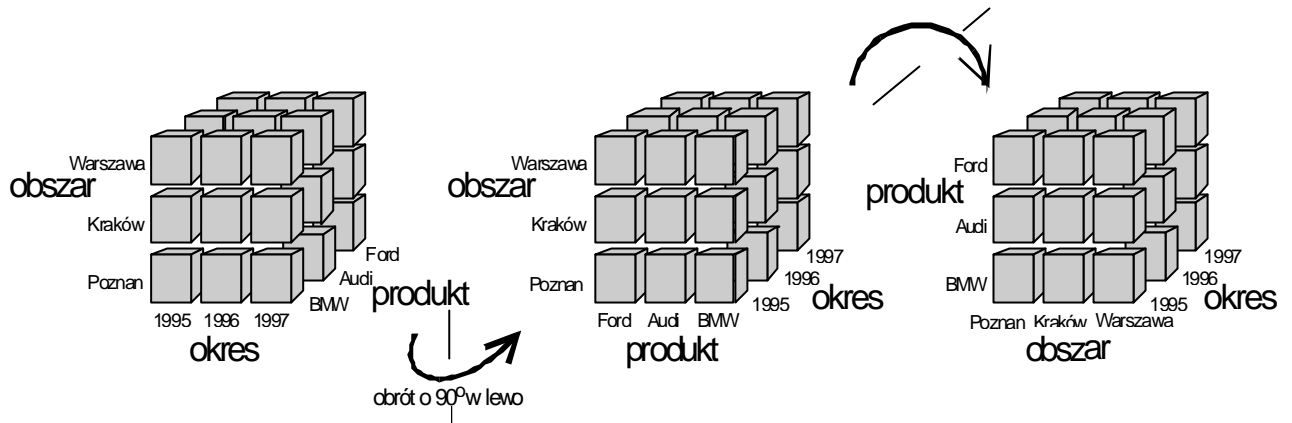
Rys.7. Operacja rozwijania hierarchii wymiaru

- *zwijanie* (ang. rolling up)

Zwijanie jest operacją odwrotną do rozwijania i polega na nawigowaniu w górę hierarchii danego wymiaru. Dzięki tej operacji można przeprowadzać analizę danych zagregowanych na wyższym poziomie hierarchii wymiarów.

- *obracanie* (ang. rotating)

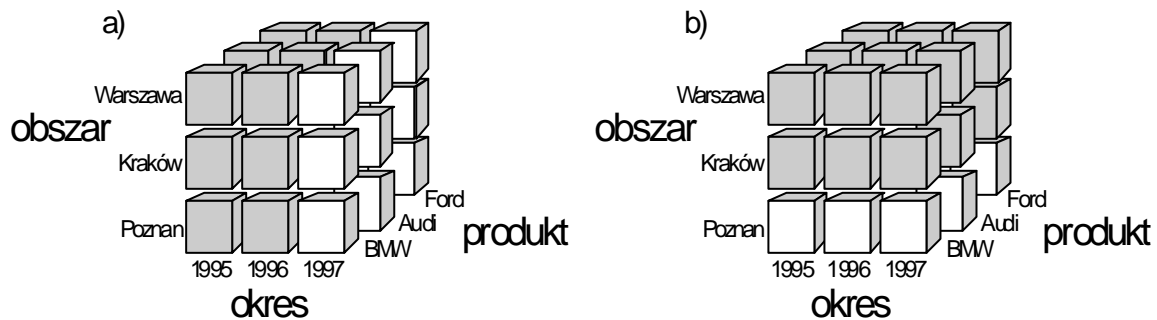
Operacja obracania (zob. Rys. 8) umożliwia prezentowanie danych w różnych układach. Celem jej jest zwiększenie czytelności analizowanych informacji.



Rys.8. Operacja obracania

- *wycinanie* (ang. slicing and dicing)

Operacja ta umożliwia zawężenie analizowanych danych do wybranych wymiarów, a w ramach każdego z wymiarów – zawężenie analizy do konkretnych jego wartości. Przykładowo, dyrektor do spraw marketingu będzie zainteresowany wielkością sprzedaży wszystkich produktów, we wszystkich miastach kraju, w roku bieżącym (zob. Rys. 9a). Natomiast kierownika oddziału firmy w Poznaniu będzie interesowała wielkość sprzedaży wszystkich produktów, w ciągu całego okresu działalności (zob. Rys. 9b).



Rys.9. Wycinanie danych w różnych wymiarach

- *obliczanie rankingu* (ang. ranking)

Operacja ta umożliwia uporządkowanie informacji w danym wymiarze, zgodnie z wartościami wybranych miar (w kolejności malejącej lub narastającej). Przykładowo, w wymiarze roku 1997 można uporządkować marki samochodów zgodnie z narastającym porządkiem liczby sprzedanych egzemplarzy.

6 Efektywność systemu hurtowni danych

Aby systemy relacyjnych baz danych mogły efektywnie przetwarzać zapytania typu OLAP i zarządzać danymi rzędu terabajtów muszą posiadać nowe właściwości zwiększające ich efektywność. Nowymi cechami systemów relacyjnych są m.in.: przetwarzanie równoległe, parcelacja danych i nowe techniki indeksowania.

6.1 Przetwarzanie równoległe

Przetwarzanie równoległe (ang. parallel processing) polega na rozbiciu złożonych operacji na mniejsze, które następnie są wykonywane równoległe, np. na wielu procesorach lub komputerach. W efekcie, czas wykonania całej operacji jest krótszy. W przypadku hurtowni danych, najczęściej równoległe przetwarzają się zapytania, sortuje dane, wykonuje operacje odczytu i zapisu na dysk, buduje tablice i indeksy oraz wczytuje dane do hurtowni.

Przetwarzanie równoległe wspierają m.in. systemy zarządzania bazami danych: *Oracle7* i *Oracle8* (Oracle Corporation), *DB2* (IBM), *OnLine Extended Parallel Server*, *OnLine Dynamic Server* (Informix), *Red Brick Warehouse* (Red Brick), *Sybase IQ* (Sybase) [3].

6.2 Parcelacja danych

Parcelacja danych (ang. data partitioning) umożliwia automatyczne rozpraszanie danych (pochodzących z jednej lub wielu relacji) na wiele dysków, znajdujących się w tym samym lub wielu węzłach (komputerach) sieci. Dzięki podziałowi dużej relacji na mniejsze:

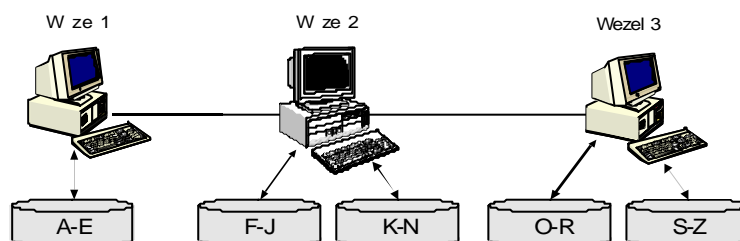
- bardzo kosztowne operacje wejścia/wyjścia, tj. dostępu do dysków mogą być wykonywane równoległe,
- równoważone jest obciążenie dysków,
- polecenia SQL mogą być wykonywane równoległe, np. tworzenie relacji i indeksów, wykonywanie zapytań,
- wzrasta bezpieczeństwo danych w przypadku awarii sprzętu,

- wzrasta szybkość tworzenia kopii zapasowych bazy i szybkość odtwarzania danych po awarii.

W komercyjnych SZBD stosuje się cztery podstawowe techniki parcelacji danych: *round–robin* (ang. round–robin partitioning), *parcelacja bazująca na wartości* (ang. range partitioning), *parcelacja haszowa* (ang. hash partitioning), *parcelacja hybrydowa* (ang. hybrid partitioning).

Technika *round–robin* umożliwia równomierne rozproszenie danych w węzłach sieci. Przykładowo, jeśli w sieci znajdują się trzy węzły, to pierwsza krotka relacji zostanie umieszczona w węzle pierwszym, druga – w węzle drugim, trzecia krotka – w węzle trzecim, czwarta – znów w węzle pierwszym itd. Rozwiązanie to ma jednak wadę. Ponieważ dane są rozproszone w sposób przypadkowy, więc odnalezienie zadanych informacji wymaga przeszukania wszystkich węzłów.

W przypadku *parcelacji bazującej na wartości* rozmieszczenie danych w sieci zależy od wartości samych danych. Przykładowo, relacja zawierająca informacje o klientach sieci supermarketów może być podzielona zgodnie z wartością pierwszej litery nazwiska, jak na Rys.8. Węzeł pierwszy posiada jeden dysk, na którym umieszczono klientów o nazwiskach rozpoczynających się od liter A do E. Węzły drugi i trzeci posiadają po dwa dyski, na których umieszczono klientów o nazwiskach odpowiednio z zakresów F–J, K–N, O–R, S–Z.



Rys. 8. Parcelacja danych bazująca na wartości

Ten sposób rozpraszania danych jest efektywny dla zapytań wykorzystujących zakresy wartości w predykatkach selekcji, ponieważ umożliwia szybki dostęp do danych z zadanego zakresu, bez potrzeby przeszukiwania wszystkich węzłów.

W przypadku *parcelacji haszowej* dane są umieszczane w węzłach zgodnie z wartością systemowej funkcji haszowej. Argumentem wejściowym tej funkcji jest wartość atrybutu, a jej wynikiem – adres węzła, w którym zostanie umieszczona krotka. W celu odnalezienia zadanych informacji SZBD wykorzystuje tę samą funkcję haszową. Zaletą tej metody jest możliwość automatycznego umieszczania w tym samym węzle krotek pochodzących z różnych, powiązanych z sobą relacji. W ten sposób zwiększa się efektywność wykonywania operacji łączenia krotek, gdyż łączone z sobą krotki znajdują się w tym samym węzle.

Parcelacja hybrydowa umożliwia dwustopniowe rozpraszanie danych. W kroku pierwszym dane są umieszczane w poszczególnych węzłach za pomocą parcelacji haszowej. W kroku drugim dane są umieszczane na poszczególnych dyskach danego węzła, za pomocą parcelacji bazującej na wartości. Dzięki tej technice wzrasta równomierność rozproszenia danych i obciążenia węzłów.

6.3 Indeksowanie danych

Indeksowanie polega na łączeniu wartości atrybutów kluczowych z adresami fizycznych bloków dyskowych w celu przyspieszenia dostępu do danych. Dla bardzo złożonych zapytań typu OLAP, operujących na ogromnej liczbie danych, standardowe indeksy w postaci B–drzew okazują się nieefektywne ponieważ

- nie zapewniają wystarczająco szybkiego dostępu do danych,
- ich rozmiar jest zbyt duży, przez co wzrastają koszty ich przetwarzania, przechowywania i utrzymywania.

Dla tej klasy zastosowań stosuje się więc nowe struktury indeksów, tj. *indeksy bitmapowe* (ang. bit–mapped indexes) i *indeksy połączeniowe* (ang. join indexes) [1], [10], [11], [12], [17].

Idea *indeksów bitmapowych* jest wykorzystanie pojedynczych bitów do zapamiętania informacji o tym, że dana wartość atrybutu występuje w określonej krotce relacji. Dla każdej unikalnej wartości atrybutu jest przechowywana tablica bitów, zwana *mapa bitowa*. Każdy bit mapy odpowiada jednej krotce relacji R – bit pierwszy odpowiada pierwszej krotce relacji R , bit drugi – drugiej krotce itp. Dla mapy $A='w'$ bit n przyjmuje wartość jeden, jeśli atrybut A krotki o numerze n przyjmuje wartość 'w'. W przeciwnym przypadku bit n przyjmuje wartość zero. Liczba bitów mapy bitowej odpowiada liczbie krotek relacji R .

Indeks bitmapowy jest zbiorem map bitowych dla wszystkich unikalnych wartości danego atrybutu. Indeks tego typu może również posiadać strukturę B–drzewa, w którego liściach zamiast adresów rekordów są przechowywane mapy bitowe.

Powyższa koncepcja zostanie zilustrowana przykładem. Tabela 1 przedstawia fragment relacji *Sprzedaz* przechowującej informacje o sprzedaży samochodów. Dla atrybutu *kolor* zdefiniowano indeks bitmapowy składający się z dwóch map bitowych. Pierwsza z nich opisuje te krotki relacji, dla których atrybut *kolor* przyjmuje wartość 'zielony', a druga – te krotki, dla których *kolor*='niebieski'. Ponieważ atrybut ten przyjmuje dwie różne wartości, więc indeks zawiera dwie mapy bitowe. Bit pierwszy mapy *kolor*='zielony' przyjmuje wartość 1, co oznacza, że atrybut *kolor* pierwszej krotki przyjmuje wartość 'zielony'. Natomiast bit drugi mapy przyjmuje wartość 0 ponieważ wartością atrybutu *kolor* krotki drugiej nie jest 'zielony'.

Sprzedaz		
klientID	marka	Kolor
1010	Fiat	zielony
1020	BMW	niebieski
1030	Fiat	zielony
1040	Audi	zielony
1050	Volvo	zielony
1060	Fiat	niebieski
1070	Ford	niebieski
1080	Opel	zielony
1090	Opel	niebieski
1100	Ford	zielony

Tabela 1. Przykładowa relacja
Sprzedaz

kolor	
zielony	niebieski
1	0
0	1
1	0
1	0
1	0
0	1
0	1
1	0
0	1
1	0

Tabela 2. Indeks bitmapowy
dla atrybutu *kolor*

Podstawowa zaleta indeksów bitmapowych jest ich mały rozmiar dla atrybutów o wąskiej dziedzinie wartości. Jako przykład rozważmy relację R zawierającą milion krotek. Przyjmijmy, że na atrybucie A przyjmującym cztery różne wartości utworzono indeks bitmapowy. Indeks ten składa się z czterech map bitowych o rozmiarze miliona bitów każda, co daje rozmiar każdej z map w przybliżeniu 125kB. Łączny rozmiar indeksu bitmapowego wynosi zatem w przybliżeniu 500kB (4 x 125kB).

Wyznaczmy teraz rozmiar tradycyjnego indeksu w postaci B–drzewa zbudowanego na atrybucie A , przyjmując następujące założenia: adresy krotek są czterobajtowe, liście zawierają skompresowane listy adresów krotek. Wówczas, indeks przechowuje milion adresów krotek, a więc jego rozmiar wynosi w przybliżeniu 4MB. Jest więc osiem razy większy od odpowiadającego mu indeksu bitmapowego.

Rozważmy jednak tę samą relację R , której atrybut A przyjmuje 64 różne wartości. W tym przypadku indeks bitmapowy będzie miał rozmiar 8MB (64 x 125kB). Będzie więc większy niż odpowiadający mu indeks w postaci B–drzewa. W celu zmniejszenia rozmiarów indeksów bitmapowych stosuje się ich automatyczną kompresję, m.in. w systemach *Oracle7*, *Oracle8* i *Sybase IQ*.

Indeksy bitmapowe są bardziej efektywne od indeksów w postaci B–drzewa tylko dla określonej klasy zapytań kierowanych do bazy danych. Są to zapytania wykorzystujące dużą liczbę predykatów warunkowych oraz zapytania wykorzystujące funkcje COUNT. Większa efektywność tych indeksów wynika z:

- dużej szybkości przetwarzania map bitowych za pomocą operatorów AND, OR i NOT – dla popularnych procesorów 64-bitowych, w jednym cyklu są przetwarzane 64 bity mapy.
- małego rozmiaru indeksów – indeksy takie zdefiniowane na atrybutach o wąskiej dziedzinie są znacznie mniejsze od indeksów w postaci B–drzewa, dzięki czemu mogą być przechowywane w pamięci operacyjnej.

- Możliwości wykonywania operacji logicznych i funkcji COUNT bezpośrednio na indeksach bitmapowych (znajdujących się w pamięci operacyjnej), a nie na samych krotkach. W rezultacie, system nie wykonuje kosztownych odczytów danych z dysku w czasie przetwarzania polecenia. Dane są odczytywane dopiero po wykonaniu wszystkich operacji logicznych na indeksach. Ich wynikiem jest mapa bitowa opisująca tylko te krotki, które spełniają warunki selekcji.

Indeksy bitmapowe wykazują jednak mniejszą efektywność dla zapytań wyszukiwujących dane z zadanych zakresów (operatory $>$, $<$ itp.), ponieważ dla takiej klasy zapytań należy wykonać wiele operacji sum logicznych na mapach bitowych dla wszystkich wartości z danego zakresu.

Rozmiar indeksu bitmapowego zależy od krotności atrybutu, dla którego został zdefiniowany. Dla danych o niskiej krotności indeks bitmapowy jest mały, w porównaniu z indeksem w postaci B-drzewa. Natomiast dla danych o wysokiej krotności indeks bitmapowy zawiera wiele map bitowych o niskiej gęstości, co pociąga za sobą spadek efektywności ich przetwarzania i zwiększa koszty składowania.

Utrzymywanie indeksów bitmapowych w czasie pracy systemu może być kosztowne, ponieważ:

- wstawienie pojedynczej krotki do relacji wymaga uaktualnienia wszystkich map bitowych zdefiniowanych dla tej relacji,
- w przypadku usuwania krotek należy utrzymywać dodatkową mapę bitową dla każdej relacji, której krotki są usuwane.

Innym rodzajem indeksu zwiększającym szybkość przetwarzania danych jest tzw. *indeks polaczeniowy* (ang. join index). Indeks tego typu łączy z sobą krotki z różnych relacji posiadające tę samą wartość atrybutu polaczeniowego, jest więc strukturą zawierającą zmaterializowane połączenie wielu relacji. Indeks taki posiada strukturę B-drzewa zbudowanego na atrybucie polaczeniowym relacji (bądź na wielu takich atrybutach). Liście indeksu zawierają wspólne wartości atrybutu polaczeniowego relacji wraz z listami adresów krotek w każdej z łączonych relacji.

Odmiana tego indeksu jest tzw. *bitmapowy indeks polaczeniowy* (ang. bit-mapped join index), który różni się od powyższego tym, że w liściach zamiast adresów krotek znajdują się mapy bitowe opisujące krotki łączonych relacji.

Oceniając przydatność indeksów bitmapowych w magazynach danych należy wziąć pod uwagę trzy następujące aspekty: zysk czasowy ze stosowania indeksów, rozmiar indeksów i łatwość utrzymywania struktury indeksów w trakcie pracy systemu.

Śród istniejących na rynku serwerów baz danych mechanizmy tworzenia i zarządzania indeksami bitmapowymi posiadają: *Oracle7* i *Oracle8* (Oracle), *Sybase IQ* (Sybase), *OnLine Dynamic*

Server, OnLine Extended Parallel Server i OnLine Workgroup Server (Informix) i Red Brick Warehouse (Red Brick Systems) [1], [6], [7], [12], [13], [14], [15]

7 Podsumowanie

W architekturze magazynów danych dane pochodzące z wielu heterogenicznych źródeł są integrowane do globalnego schematu magazynu danych. Dane są przechowywane w magazynie, który z punktu widzenia użytkownika jest zwyczajną bazą danych. Po załadowaniu danych do magazynu użytkownicy mogą wydawać zapytania do magazynu (analogicznie do zwyczajnej bazy danych), ale, generalnie, nie wolno im modyfikować danych znajdujących się w magazynie, ponieważ takie modyfikacje nie znajdowałyby odzwierciedlenia w zmianach danych źródłowych i mogłyby prowadzić do niespójności magazynu. Istnieją trzy główne podejścia do zagadnienia aktualności magazynu danych [4]:

- Magazyn jest okresowo rekonstruowany z danych źródłowych. Rekonstrukcja przebiega najczęściej w nocy (gdy z magazynu nie korzystają żadne aplikacje) i wymaga wyłączenia magazynu. Wadą tej metody jest czasochłonność oraz fakt, że dane w magazynie nie odzwierciedlają aktualnego stanu danych źródłowych.
- Magazyn danych jest okresowo aktualizowany i zmiany, które zaszły w źródłach, są propagowane do magazynu. Takie podejście wymaga uaktualnienia znacznie mniejszej ilości danych. Z drugiej strony inkrementalna aktualizacja jest znacznie trudniejsza a algorytmy wykrywające zmiany w danych źródłowych są bardziej skomplikowane.
- Magazyn danych jest aktualizowany natychmiast po wystąpieniu zmiany w źródle. Ze względu na ilość komunikacji i przetwarzania związanego z takim rozwiązaniem, znajduje ono zastosowanie tylko w małych magazynach danych, w których dane źródłowe zmieniają się stosunkowo powoli.

Magazyny danych odgrywają szczególną rolę w aplikacjach OLAP. Po pierwsze, magazyn danych jest potrzebny do scentralizowania i zorganizowania danych w taki sposób, by efektywnie wspierać aplikacje OLAP. Dane potrzebne do analizy mogą być rozproszone w wielu bazach danych przedsiębiorstwa. Poza tym zapytania OLAP są złożone i wykorzystują duże wolumeny danych, stąd ich wykonanie w systemach zorientowanych na przetwarzanie transakcji (czyli systemach OLTP) jest z reguły nieefektywne.

Aplikacje OLAP wykorzystują wielowymiarowy model danych. W modelu tym dane prezentowane są jako komórki w wielowymiarowej przestrzeni i nazywane są miarami. Istnieją dwie podstawowe metody implementacji takiego modelu danych:

- ROLAP – dane są przechowywane w relacyjnych tabelach ułożonych w schemat gwiazdy. W centrum takiej gwiazdy znajduje się relacja faktów, przechowująca nieprzetworzone informacje o miarach, zaś wymiary przechowywane są w relacjach wymiarów.
- MOLAP – dane wielowymiarowe przechowywane są w wyspecjalizowanych strukturach nazywanych kostkami danych. Dane przechowywane w kostce danych mogą być wstępnie przetworzone, oprócz tego dla takiej kostki mogą być zdefiniowane specjalne operatory, np. rozwijanie i zwijanie wymiarów, obracanie, wycinanie, itp.

W dziedzinie magazynów danych dużego znaczenia nabiera zagadnienie efektywności przetwarzania złożonych i czasochłonnych zapytań, które przetwarzają duże ilości danych. W celu poprawy efektywności działania magazynów danych stosuje się wiele technik, których celem jest poprawa owej efektywności. Należą do nich m.in. przetwarzanie równoległe, parcelacja danych oraz nowe metody indeksowania, np. omówione w niniejszym artykule indeksy bitmapowe.

8 Bibliografia

- [1] Abbey M., Corey M. J., *Oracle8 – A Beginner's Guide*, Osborne McGraw–Hill, 1997, ISBN 0-07-882393-5
- [2] Colliat G., *OLAP, Relational, and Multidimensional Database Systems*, SIGMOD Record, Vol.25, No.3, wrzesień 1996
- [3] Data Management Strategies, Cutter Information Corp., <http://www.cutter.com/itgroup/>
- [4] Garcia-Molina, H., Ullman, J.D., Widom, J., *Database System Implementation*, Prentice Hall, 2000
- [5] Hurwicz, M., *Take Your Data To The Cleaners*, BYTE, styczeń 1997
- [6] *INFORMIX–OnLine Extended Parallel Server for Loosely Coupled Cluster and Massively Parallel Processing Architectures*, materiały informacyjne firmy Informix, 1997
- [7] *INFORMIX–Data Warehouse Servers. The Evolution of Integrated Relational OLAP*, Informix White Paper, 1996,
<http://www.informix.com/informix/corpinfo/zines/whitpprs/dwserver/dwserver.htm>
- [8] *Designing the Data Warehouse on Relational Databases*,
<http://www.informix.com/informix/corpinfo/zines/whitpprs/stg/metacube.htm>
- [9] Jarke, M., Lenzerini, M., Vassiliou, Y., Vassiliadis, P., *Fundamentals of Data Warehouses*, Springer-Verlag Berlin Heidelberg, 2000
- [10] O'Neil P., Graefe G., *Multi-Table Joins Through Bitmapmed Join Indices*, SIGMOD Record, Vol.24, No.3, wrzesień 1995
- [11] O'Neil P., Quass D., *Improved Performance with Variant Indexes*, materiały konferencyjne SIGMOD, Tuscon, Arizona, USA, 1997
- [12] *Bitmapmed Indexes in Oracle7*, an Oracle White Paper, 1995
- [13] *A Red Brick Systems White Paper*, http://www.redbrick.com/rbs/whitepapers/datawah_wp.html
- [14] Sarawagi S., *Indexing OLAP data*, Data Engineering Bulletin 20(1), 1997
- [15] *Optimizing Interactive Performace for the Data Warehouse*,
<http://www.novasys.com.cy/sysprod06.thm>
- [16] Widom J., *Research Problems in Data Warehousing*, materiały konferencyjne, 4–th International Conference on Information and Knowledge Management (DIKM), 1995
- [17] Wrembel R., *Nowe struktury indeksów dla magazynów danych*, materiały konferencyjne III Jesiennej Szkoły PLOUG, Zakopane, 1997