

Aktywne hurtownie danych

Mikolaj Morzy

Instytut Informatyki Politechniki Poznańskiej
ul. Piotrowo 3A, 60-965 Poznan
Mikolaj.Morzy@cs.put.poznan.pl

Abstract. Tradycyjne hurtownie danych stały się nieodłącznym składnikiem infrastruktury informatycznej wielu przedsiębiorstw. Wraz z gwałtownie następującym postępem technicznym pojawiają się jednak wyzwania, którym tradycyjna architektura hurtowni danych nie może sprostać. W szczególności, tradycyjne hurtownie danych wspierają działalność przedsiębiorstwa na szczeblu strategicznym, lecz nie oferują taktycznego wsparcia podejmowania decyzji na szczeblu operacyjnym. Aby sprostać tym wymaganiom, zaproponowano nowy model informatycznej infrastruktury przedsiębiorstwa, zwany aktywną hurtownią danych. Aktywna hurtownia danych umożliwia współdzielenie danych składowanych w operacyjnych bazach danych z danymi przechowywanymi w hurtowniach danych i stanowi istotny komponent umożliwiający integrację aplikacji działających w przedsiębiorstwie. W niniejszym artykule przedstawiamy pojęcie aktywnej hurtowni danych i wskazujemy na czynniki, które decydują o powodzeniu wdrożenia takiego rozwiązania. Omawiamy również korzyści, jakie mogą zostać odniesione przez przedsiębiorstwo w efekcie wdrożenia aktywnej hurtowni danych. Dyskusję ilustrujemy przykładami wdrożeń, które zakończyły się sukcesem. Artykuł zamyka dyskusja na temat integracji aplikacji przedsiębiorstwa i roli, jaką w tym procesie odgrywa aktywna hurtownia danych.

Wprowadzenie

Model przetwarzania analitycznego on-line

Sposób, w jaki użytkownik korzysta z systemu komputerowego, nazywamy modelem przetwarzania. W dziedzinie baz danych najbardziej rozpowszechnionym modelem przetwarzania jest *przetwarzanie transakcji w trybie on-line* (ang. on-line transaction processing, OLTP). Model ten charakteryzuje się założeniem, że w systemie wykonują się współbieżnie krótkie transakcje, z których każda dokonuje odczytu lub zapisu niewielkiej ilości danych, zazwyczaj kilku rekordów. Model OLTP zakłada również, że równolegle wykonują się setki lub tysiące transakcji o podobnej charakterystyce. W systemach implementujących taki model przetwarzania kluczowe znaczenie mają zagadnienia takie jak zapewnienie spójności danych, integralność danych, maksymalizacja współbieżności przetwarzania, zapewnienie odpowiedniego poziomu izolacji między transakcjami, czy udostępnienie wydajnych mechanizmów wycofywania lub odtwarzania transakcji, które uległy uszkodzeniu. Model OLTP dobrze charakteryzuje przetwarzanie typowe dla systemów służących do obsługi bieżącej działalności przedsiębiorstwa, takich jak systemy księgowe, finansowe, kadrowe, czy magazynowe.

Współczesne przedsiębiorstwa wymagają jednak bardziej zaawansowanych narzędzi informatycznych aby radzić sobie w obliczu ostrej konkurencji. Jednymi z najbardziej użytecznych narzędzi są systemy wspomagania decyzji. Systemy te umożliwiają analizę i wykorzystanie ogromnych ilości danych gromadzonych przez przedsiębiorstwa. Większość przedsiębiorstw, organizacji, instytucji publicznych, czy ośrodków naukowych gromadzi ogromne ilości danych. Dane te, przechowywane w heterogenicznych bazach danych, arkuszach kalkulacyjnych lub dokumentach tekstowych, są dodatkowo rozproszone geograficznie i niespójne ze względu na lokalizację narodową. Mimo, że ich integracja jest zadaniem trudnym, niemniej jednak jest opłacalna. Wynika to z faktu, że zgromadzone archiwa zawierają w sobie bardzo dużą ilość wartościowej wiedzy, która może być wykorzystana do wspierania działalności przedsiębiorstwa, w szczególności do wspomagania decyzji. Informacje zawarte w zgromadzonych danych opisują trendy, anomalie, regularności, korelacje i użyteczne wzorce, np. wzorce zachowań klientów, wzorce okresowych wahań sprzedaży, itp. Aby umożliwić wykorzystanie danych gromadzonych przez przedsiębiorstwo, zaproponowano nowy model przetwarzania, nazwany *przetwarzaniem analitycznym on-line* (ang. on-line analytical processing) i opracowano nową architekturę systemu baz danych, nazwaną hurtownią danych.

Model przetwarzania analitycznego on-line zakłada, że w systemie wykonuje się współbieżnie niewiele transakcji (zazwyczaj kilka do kilkunastu), przy czym każda transakcja posiada dwie charakterystyczne cechy: dokonuje tylko i wyłącznie odczytów danych i przetwarza jednorazowo bardzo duży wolumen danych rzędu tysięcy lub milionów rekordów. Wydajność przetwarzania mierzy się nie, jak w przypadku modelu OLTP, liczbą transakcji w jednostce czasu, lecz czasem wykonania pojedynczej transakcji. Ponieważ transakcje dokonują

tylko i wyłącznie odczytów danych, stąd zagadnienia kluczowe dla modelu OLTP, np. synchronizacja współbieżnego wykonania zbioru transakcji, w modelu OLAP stają się nieistotne. Z drugiej strony, dla efektywności przetwarzania ważne stają się optymalizacja zapytań, efektywne struktury dostępu do danych, wstępna agregacja i materializacja danych.

Hurtownia danych

Model przetwarzania OLAP stał się podstawą hurtowni danych. Hurtownia danych to „(...) zorientowana tematycznie, zintegrowana, zmienna w czasie i nieulotna kolekcja danych wspierająca proces wspomaganie decyzji” (W.H.Inmon). Hurtownia danych jest zorientowana tematycznie, co oznacza, że dane przechowywane w hurtowni danych stanowią wyczerpujący opis pewnego fragmentu rzeczywistości, w której działa dane przedsiębiorstwo. Zazwyczaj perspektywa danych przechowywanych w hurtowni danych różni się znacząco od perspektywy danych przechowywanych w operacyjnej bazie danych. Dla przykładu, dane w systemie transakcyjnym banku przedstawiają działalność banku jako sekwencje transakcji, podczas gdy hurtownia danych banku jest zorientowana na klientów i ich konta, lokaty, kredyty, itp. Integracja hurtowni danych polega na tym, że dane przechowywane w hurtowni pochodzą z wielu potencjalnie heterogenicznych źródeł danych. Integracja danych z wszystkich gałęzi działania przedsiębiorstwa jest kluczem do osiągnięcia wysokiej jakości hurtowni danych. Dane w przedsiębiorstwie są przechowywane w wielu izolowanych systemach, działających na różnych platformach i stworzonych w różnych technologiach. Integracja polega na ich scaleniu, ujednoczeniu formatów i nazewnictwa, konwersji danych, itp. Zmienność w czasie spowodowana jest tym, że hurtownia danych zawiera dane historyczne. Hurtownie danych można postrzegać jako zbiór migawek, z których każda przedstawia globalny stan informacji obecnej wewnątrz przedsiębiorstwa w danym momencie. Horyzont czasowy jest bardzo istotną składową hurtowni danych, umożliwiającą analizę trendów i wzorców zmiennych w czasie oraz ewolucję rzeczywistości, w której działa dane przedsiębiorstwo. Wreszcie nieulotność danych zgromadzonych w hurtowni danych powoduje, że dane, które zostały wprowadzone do hurtowni danych nigdy nie ulegają żadnym modyfikacjom. Najczęściej jedynymi dozwolonymi modyfikacjami zawartości hurtowni danych jest odświeżenie hurtowni (wprowadzenie nowych danych z systemów źródłowych) i archiwizacja [1,4].

Hurtownie danych są narzędziem skierowanym przede wszystkim do analityków, strategów i decydentów. Umożliwiają im analizę dużych ilości danych i określanie trendów, potwierdzanie lub obalanie hipotez, identyfikację zasobów (klientów, produktów lub usług) przynoszących wysoki dochód lub wysokie straty. Analiza zawartości hurtowni danych najczęściej odbywa się na podstawie zapytań formułowanych przez użytkownika, gdzie poszczególne zapytania analizują dane atomowe zgromadzone w tzw. *tabeli faktów* (ang. fact table) w kontekście powiązanych z faktami wymiarów. Fakty reprezentują podstawowe „cegielki” składające się na analizowaną rzeczywistość (np. pojedynczy fakt może reprezentować sprzedaż produktu, podpisanie umowy z klientem, lub wykonanie połączenia telefonicznego) i zazwyczaj przechowują mierzalną informację o danym fakcie, zwana *miarą* (ang. measure). Przykładem miary może być cena, liczba sprzedanych produktów, lub czas połączenia telefonicznego. Wymiary, przechowywane w *tabelach wymiarów* (ang. dimension tables) służą do opisywania faktów. Przykładami wymiarów są: czas, klient, produkt, czy abonament telefoniczny. Wymiary mogą tworzyć hierarchie umożliwiające na analizę danych z podziałem na miasta, powiaty, województwa (hierarchia wymiaru opisującego lokalizację) lub z podziałem na miesiące, kwartały i lata (hierarchia wymiaru opisującego czas). Zapytania kierowane przez użytkowników do hurtowni danych dokonują wyliczenia złożonych agregatów wykorzystując do tego zdefiniowane wymiary. Przykładem zapytania analitycznego może być zadanie wyliczenia średniej sprzedaży danej grupy produktów z podziałem na lokalizację sklepu, kategorie klientów i miesiące. Zapytania takie są formułowane przez analityków, menedżerów i decydentów w celu budowania strategicznej wizji działalności przedsiębiorstwa. Warto tutaj zauważyć, że grono użytkowników tradycyjnej hurtowni danych jest bardzo ograniczone. Co więcej, wzorce odkryte w hurtowni danych mogłyby wspierać podejmowanie decyzji taktycznych w przedsiębiorstwie, jeśli istniałaby możliwość udostępnienia tych danych użytkownikom końcowym (sprzedawcom, urzędnikom, dostawcom, itp.). Współcześnie dostępne hurtownie danych nie posiadają takich możliwości. Właśnie ten fakt stał się źródłem krytyki wielu specjalistów i spowodował opracowanie nowej architektury aktywnych hurtowni danych.

Operacyjna składnica danych

Podstawowym problemem informatycznym z jakim borykają się współczesne przedsiębiorstwa jest brak zintegrowanego i aktualnego obrazu danych wykorzystywanych przez dane przedsiębiorstwo. Jak wskazano wcześniej, tradycyjna hurtownia danych tylko częściowo rozwiązuje ten problem. Po pierwsze, dane przechowywane w hurtowni nigdy nie są aktualne i w zależności od procedury odświeżania hurtowni mogą

reprezentować migawkę sprzed wielu dni lub tygodni. Po drugie, dane w hurtowni często są już przetworzone do postaci wymaganej przez narzędzia analityczne i raportujące. Taka transformacja może za sobą pociągać uogólnienie danych i utratę informacji o danych atomowych. Konieczność posiadania zintegrowanych i aktualnych danych rodzi potrzebę budowania *operacyjnej składnicy danych* (ang. operational data store, ODS).

Konieczność wyodrębnienia operacyjnej składnicy danych

Hurtownia danych przedsiębiorstwa jest częścią modułu *Inteligencji Biznesowej* (ang. Business Intelligence). Oprócz samej hurtowni danych na modul Inteligencji Biznesowej składają się aplikacje analityczne i raportujące, oraz *tematyczne hurtownie danych* (ang. data marts). Aplikacje analityczne komunikują się z hurtownią danych i umożliwiają analitykom i decyzyjcom wgląd w statystyki, podsumowania, informacje demograficzne, prognozy, itp. Informacje uzyskane za pomocą aplikacji analitycznych służą do tworzenia raportów, zestawień i podsumowań. Problem polega na tym, że aktualna infrastruktura nie umożliwia wykorzystania wiedzy pochodzącej z analizy hurtowni danych do podjęcia decyzji taktycznych, czyli takich, które są podejmowane w aplikacjach operacyjnych. Aplikacje operacyjne stanowią komponenty *modułu Zarządzania* (ang. Business Management). Komponenty te umożliwiają przedsiębiorstwu wykonywanie swojej misji na podstawie wiarygodnych danych i prognoz. Jakość usług świadczonych przez przedsiębiorstwo, a zatem jakość procesów wykonywanych za pomocą aplikacji operacyjnych, w dużej mierze zależy od jakości danych dostarczanych do aplikacji operacyjnych, a ta z kolei zależy od stopnia integracji modułu Zarządzania z modulem Inteligencji Biznesowej. Modul Zarządzania jest wykorzystywany przez personel przedsiębiorstwa mający bezpośredni styk z klientem, aplikacje wchodzące w skład tego modułu muszą mieć do dyspozycji pełną, zintegrowaną, aktualną i wzbogaconą wiedzę o kliencie. Należy zauważyć, że informacje dostarczane do aplikacji operacyjnych muszą stanowić połączenie aktualnych i szczegółowych danych tradycyjnie przechowywanych w operacyjnych bazach danych z wynikami analiz (np. danymi demograficznymi) powstałych w hurtowni danych. Architektura umożliwiająca połączenie aplikacji operacyjnych z aplikacjami analitycznymi i stosująca zaawansowaną integrację danych zwana jest *korporacyjną fabryką informacji* (ang. Corporate Information Factory, CIF) i stanowi bardzo obiecującą propozycję, która w przyszłości bez wątpienia będzie odgrywała istotną rolę w infrastrukturze informatycznej nowoczesnych przedsiębiorstw.

Operacyjna składnica danych, podobnie jak hurtownia danych, jest zorientowana tematycznie. Oznacza to, że dane przechowywane w operacyjnej składnicy danych przedstawiają zintegrowany fragment rzeczywistości w której działa przedsiębiorstwo i udostępniają wgląd do danego wycinka danych z perspektywy całego przedsiębiorstwa, a nie tylko wybranego działu (jak dzieje się w przypadku zwykłej operacyjnej bazy danych). Dla przykładu, operacyjna składnica danych zorientowana na klienta może zawierać, poza szczegółowymi danymi o kliencie, dane o wszystkich ostatnich interakcjach klienta z przedsiębiorstwem (listy zakupionych produktów, stan rozliczeń z klientem, informacje o kontaktach telefonicznych klienta z działem obsługi, itp.). Dane przechowywane w operacyjnej składnicy danych są też zintegrowane i stanowią spójny katalog danych. Kontynuując przykład, operacyjna składnica danych zawiera wszystkie informacje o danym kliencie zebrane w całym przedsiębiorstwie. Dzięki temu spójny katalog danych o kliencie jest wykorzystywany we wszystkich punktach styczności klienta z przedsiębiorstwem i może być wykorzystywany przez cały personel który wchodzi w interakcje z danym klientem. W przeciwieństwie do hurtowni danych operacyjna składnica danych nie zawiera horyzontu czasowego. Dane przechowywane w OSD są zawsze najświeższe, podobnie jak w tradycyjnej operacyjnej bazie danych. Oczywiście, nic nie stoi na przeszkodzie aby w ramach spójnego katalogu integrującego dane o konkretnym kliencie zawrzeć również pewną ilość danych historycznych (np. poprzedni adres, poprzedni telefon kontaktowy, transakcje z ostatniego tygodnia) jeśli istnieją aplikacje operacyjne wykorzystujące takie dane. Generalnie jednak rzecz ujmując, operacyjna składnica danych jest ahistoryczna. Ponieważ dane prezentowane w operacyjnej składnicy danych są nieustannie aktualizowane, są one również ulotne (to kolejna istotna różnica między OSD a hurtownią danych). Wszystkie aktualizacje dokonywane w operacyjnych bazach danych muszą być jak najszybciej propagowane do OSD, aby zapewnić wszystkim użytkownikom OSD aktualny widok danych. Wreszcie OSD zawiera dane szczegółowe, bez jakichkolwiek wstępnie wyliczonych agregatów i podsumowań. W przeciwieństwie do hurtowni danych, szczegółowe dane przechowywane w OSD mają charakter dynamiczny i często ulegają modyfikacjom.

Odświeżanie operacyjnej składnicy danych

Jedną z najważniejszych różnic występujących między operacyjną składnicą danych a tradycyjną hurtownią danych jest częstotliwość odświeżania. Tradycyjna hurtownia danych, odświeżana za pomocą metod ETL (Extract-Transform-Load) lub ELT (Extract-Load-Transform), jest najczęściej odświeżana okresowo, w momentach bezczynności. W zależności od domeny zastosowania hurtownie danych są odświeżane w cyklach

dziennych, tygodniowych lub miesięcznych. Operacyjne składnice danych można sklasyfikować ze względu na częstotliwość odświeżania w następujący sposób:

- Klasa I: dane wprowadzone do aplikacji operacyjnej muszą zostać propagowane do operacyjnej składnicy danych natychmiast, w przeciągu kilku sekund. Przykładem takich danych mogą być dane nowego klienta, które podaje on za pomocą formularza na stronie internetowej przedsiębiorstwa.
- Klasa II: dane mogą zostać wprowadzone do operacyjnej składnicy danych w sposób asynchroniczny, lecz czas odświeżenia nie powinien wynieść więcej niż godzinę. Przykładem takich danych są podsumowania opisujące profil zakupów klienta. Takie dane są wykorzystywane m.in. do wygenerowania automatycznych rekomendacji dla danego klienta i w takim celu wystarczający jest przybliżony profil danego klienta.
- Klasa III: dane należące do tej klasy mogą być odświeżane w cyklach dziennych. Dotyczy to w szczególności danych pochodzących z automatycznych systemów odkrywania wiedzy, np. dane o wzorcach zakupów, na podstawie których buduje się automatyczne rekomendacje, nie zmieniają się na tyle często, aby konieczne było ich nieustanne odświeżanie.
- Klasa IV: do tej klasy należą dane pochodzące z analiz przeprowadzanych w hurtowni danych. Częstotliwość odświeżania tych danych zależy od częstotliwości odświeżania hurtowni danych. Przykładami takich danych mogą być uaktualnienia statystyk opisujących grupy klientów lub zmieniające się trendy sprzedaży. Fakt, iż dane te są uaktualniane stosunkowo rzadko nie oznacza bynajmniej, że dane te są mniej istotne. Z analizy hurtowni danych może wynikać np. model opisujący prawdopodobieństwo tego, że dany klient zrezygnuje z usług przedsiębiorstwa na rzecz konkurencji. Znajomość i jakość tego modelu jest kluczowa w pracy personelu bezpośrednio kontaktującego się z klientem.

W celu dokładnego zobrazowania idei operacyjnej składnicy danych rozważmy następujący przykład. Klient loguje się do sklepu internetowego. System identyfikuje klienta i odczytuje wyniki analizy zachowań klienta. Na wynik takiej analizy składają się różne czynniki, między innymi historia nawigacji po sklepie, oglądane artykuły, dokonane poprzednio zakupy, dane demograficzne klienta, informacje z centrum obsługi nt. pytań zadawanych przez klienta, itp. Powyższe dane, zintegrowane i połączone, pozwalają przedsiębiorstwu na zbudowanie dokładnego profilu klienta obejmującego jego zainteresowania, zwyczaje oraz preferowane metody zakupu i zapłaty. Na bazie tych informacji system może wygenerować ofertę promocyjną, która będzie idealnie dostosowana do wymagań danego klienta. Istotne jest to, że cały proces analizy nie może zachodzić w czasie rzeczywistym, ponieważ system nie ma na to czasu. Od momentu przyłączenia się klienta do momentu wyświetlenia mu oferty nie może upłynąć więcej niż 5-6 sekund. Taki czas jest wysoce niewystarczający aby przeanalizować historię nawigacji, demografię klienta, historię jego zakupów, itd. Wszystkie te dane muszą już być obecne w operacyjnej składnicy danych i podczas logowania się klienta odpowiedni profil musi być tylko zlokalizowany. Na podstawie profilu system powinien maksymalnie szybko zaproponować oferowany produkt promocyjny. Przeciwnościem takiego podejścia jest metoda często wykorzystywana przez różnych sprzedawców i zwana pogardliwie „spray and pray”, a polegająca na zarzuceniu wszystkich klientów licznym zbiorem takich samych ofert i oczekiwaniu, że nikły procent klientów zainteresuje się niektórymi produktami.

Aktywna hurtownia danych

Aktywna hurtownia danych (ang. active data warehouse) to środowisko zapewniające zintegrowane składowanie spójnych danych na potrzeby zarówno strategicznego wspomaganie decyzji jak i taktycznych aplikacji operacyjnych. Innymi słowami, aktywna hurtownia danych to technologia umożliwiająca połączenie w ramach jednego produktu zarówno tradycyjnej hurtowni danych, jak i operacyjnej składnicy danych. Z racji swego dualnego charakteru aktywna hurtownia danych posiada pewne interesujące własności, które zostały podsumowane poniżej.

Aktywna hurtownia danych nie może ograniczać się do efektywnego wsparcia jednego modelu przetwarzania. Ponieważ na obciążenie aktywnej hurtowni danych składa się mieszanka złożonych zapytań analitycznych (model OLAP) i krótkich, szybkich transakcji pochodzących z aplikacji operacyjnych (model OLTP), architektura aktywnej hurtowni danych musi spełniać zupełnie nowe, czasem przeciwstawne wymagania dotyczące efektywności przetwarzania, skalowalności i niezawodności. Przykładem sprzecznych wymagań jest wymaganie dotyczące szybkości odpowiedzi. W danym momencie w systemie aktywnej hurtowni danych może znajdować się kosztowne i czasochłonne zapytanie wyliczające współczynniki demograficzne dla jakiejś grupy klientów. W tym samym czasie z aplikacji operacyjnej do systemu może zostać wprowadzone zapytanie wyszukujące najlepszą ofertę dla konkretnego klienta. Ponieważ drugie zapytanie ma wyższy priorytet (prawdopodobnie osoba wprowadzająca to zapytanie właśnie rozmawia z klientem), pierwsze zapytanie powinno zostać wstrzymane lub przełączone na niższy priorytet w celu zagwarantowania jak najszybszej odpowiedzi na drugie zapytanie. Z drugiej strony, mocne obciążenie ze strony aplikacji operacyjnych może oznaczać całkowite zablokowanie analitycznych zdolności aktywnej hurtowni danych.

Zazwyczaj rozmiar hurtowni danych jest większy niż rozmiar operacyjnej bazy danych. Ponieważ aktywna hurtownia danych łączy ze sobą dane tradycyjnie przechowywane w hurtowni danych z danymi operacyjnymi, skalowalność takiego rozwiązania staje się istotnym problemem. Architektura aktywnej hurtowni danych musi zapewniać efektywne wykonywanie współbieżnych zapytań działających na bardzo dużych wolumenach danych. Jednak najistotniejsza różnica między hurtownią danych a systemem operacyjnym jest zagadnienie dostępności. Tradycyjne hurtownie danych nie muszą być dostępne 24×7×365. Aktywna hurtownia danych musi również wspierać działalność operacyjną przedsiębiorstwa. Ze względu na gwałtowny wzrost internetowego kanału działalności przedsiębiorstw (przede wszystkim chodzi o sprzedaż, ale nie tylko), infrastruktura informatyczna leżąca u podstaw funkcjonowania aplikacji operacyjnych musi być dostępna zawsze. Oznacza to, że aktywna hurtownia danych nie może mieć żadnych przestoju, nie posiada okresów bezczynności przeznaczonych tradycyjnie na odświeżanie, oraz musi być w 100% niezawodna.

Innym czynnikiem, który istotnie różni aktywną hurtownię danych od tradycyjnej hurtowni danych jest kwestia świeżości danych. Aktywna hurtownia danych musi działać na najnowszych danych, co powoduje, że mechanizmy dostarczania danych do hurtowni danych muszą ulegać zmianie. Tradycyjny model asynchronicznej propagacji zmian do hurtowni danych nie zdaje egzaminu. W idealnym przypadku proces odświeżania jest ciągły i zapewnia całkowitą synchronizację danych w aktywnej hurtowni danych. Jak wspomniano wcześniej, wymagany stopień świeżości danych zależy od konkretnej aplikacji. Należy jednak przyjąć, że aktywna hurtownia danych wymaga w ogólności danych o bardzo wysokim stopniu świeżości [3].

Zyski wynikające z aktywnej hurtowni danych

Implementacja aktywnej hurtowni danych może zaowocować wieloma istotnymi zyskami. Architektura aktywnej hurtowni danych charakteryzuje się skalowalnością, dostępnością, oraz efektywnością. Obecność aktywnej hurtowni danych eliminuje nadmiarowość danych, ponieważ aktywna hurtownia danych dostarcza jednego spójnego obrazu wszystkich danych jakimi dysponuje przedsiębiorstwo. Takie zjawisko nazywa się *jedyną wersją prawdy* (ang. *single version of truth*) i ma niebagatelne znaczenie praktyczne. Dzięki brakowi redundancji danych przedsiębiorstwo unika podejmowania sprzecznych decyzji dotyczących tego samego klienta, nie zmusza klienta do wielokrotnego podawania tych samych danych, unika nieprzyjemnych sytuacji i nieporozumień wynikających z nadmiarowości danych. Istnienie jednej wersji prawdy w znaczący sposób zwiększa też zaufanie między klientem i przedsiębiorstwem.

W tradycyjnej hurtowni danych wyniki analiz są dostępne tylko analitykom, menedżerom i decydentom. Historia hurtowni danych wskazuje jasno, że wyposażenie decydentów w narzędzia analityczne znacząco poprawia jakość podejmowanych przez nich decyzji. Aktywna hurtownia danych pozwala na propagację tego zjawiska w dół hierarchii przedsiębiorstwa, do szeregowych pracowników. Oczywiście, zamiast udostępniać analizy i opracowania globalne, uzyskują one istotnie wzbogacone dane dotyczące aktualnie obsługiwanego klienta. Podobnie jak w przypadku decydentów wyższego szczebla, pracownicy niższych szczebli są w stanie podjąć lepsze decyzje jeśli są w posiadaniu lepszych danych.

Aktywna hurtownia danych minimalizuje również czas opóźnienia między podjęciem decyzji strategicznej i taktycznej. Wartość podjętej akcji szybko maleje wraz z czasem. Jeśli przedsiębiorstwo jest w stanie zareagować na jakies zdarzenie natychmiast po jego wystąpieniu, wartość takiej reakcji jest dużo wyższa niż wartość późniejszej reakcji. Przykładowo, zaoferowanie klientowi towaru po promocyjnej cenie w momencie, gdy klient znajduje się przy kasie sklepowej ma dużo większą wartość (i niższą cenę) niż zaoferowanie tego samego produktu w katalogu wysłanym do klienta. Wynika stąd, że oprócz szybkości wprowadzenia danych do aktywnej hurtowni danych nie mniej istotna jest szybkość uzyskania sprzężenia zwrotnego, czyli wsparcie decyzji taktycznej która odbywa się w punkcie styczności z klientem. Aktywna hurtownia danych znacząco zmniejsza opóźnienie dostarczenia decyzji taktycznej, ponieważ hurtownia danych będąca źródłem decyzji i aplikacja operacyjna będąca konsumentem decyzji współdzieli to samo środowisko.

Do istotnych zalet aktywnej hurtowni danych zaliczyć również należy automatyzację procesów związanych z działaniem aplikacji operacyjnych. Ponieważ w środowisku aktywnej hurtowni danych aplikacje operacyjne dysponują danymi dużo wyższej jakości niż ma to miejsce w tradycyjnym środowisku, część procesów może zostać w pełni zautomatyzowana. Dzięki temu aplikacje operacyjne mogą zostać wyposażone w znacznie większą ilość „inteligencji”, bardziej polegając na danych dostarczonych przez aktywną hurtownię danych niż na arbitralnych decyzjach użytkownika.

Od czego zależy sukces?

Implementacja aktywnej hurtowni danych jest trudnym i wymagającym zadaniem. Można pokusić się o wskazanie istotnych czynników, które decydują o sukcesie tego przedsięwzięcia. Pierwszym czynnikiem, który należy wziąć pod uwagę, jest zapewnienie wsparcia ze strony kadry zarządzającej. Należy zdawać sobie sprawę

z tego, że implementacja hurtowni danych oznacza daleko idące zmiany w polityce funkcjonowania przedsiębiorstwa. Takie zmiany mogą napotkać na opór, w szczególności ze strony szeregowych pracowników, których taka zmiana najbardziej dotyczy. Stąd zyskiwanie zaufania kierownictwa i przyszłych użytkowników ma kluczowe znaczenie dla powodzenia całego projektu. Drugim czynnikiem jest edukowanie użytkowników i promowanie otwartej polityki informacyjnej. W ten sposób unika się sytuacji, w której dział IT staje się wąskim gardłem ograniczającym użytkownikom dostęp do informacji zawartej w hurtowni danych. Kolejnym czynnikiem jest uproszczenie architektury. Aktywna hurtownia danych zazwyczaj jest implementowana w przedsiębiorstwach posiadających długą historię wdrożeń różnych produktów informatycznych. Aby aktywna hurtownia danych nie okazała się kolejnym, nieudanym wdrożeniem, ważne jest podejście totalne: usunięcie naleciałości i starych systemów, jeśli ich funkcjonalność może być w łatwy sposób zapewniona przez aktywną hurtownię danych. Doświadczenie wskazuje, że ten krok również często napotyka na znaczny opór użytkowników, którzy życzą sobie, aby nowo wdrażany produkt w niczym nie zmieniał procedur, których się nauczyli i do których zdążyli się przyzwyczaić. Ostatni czynnik nie jest charakterystyczny tylko dla aktywnych hurtowni danych, ale odnosi się do każdego nowego produktu informatycznego. Ze względu na bardzo dynamiczny charakter dziedziny jaka jest informatyka wybór stabilnej platformy i wiarygodnego dostawcy oprogramowania nabiera szczególnego znaczenia w przypadku dużego i kosztownego wdrożenia, a do takich bez wątpienia należy aktywacja hurtowni danych przedsiębiorstwa.

Wyzwania technologiczne

Ponieważ, jak już wcześniej wielokrotnie wspomiano, aktywne hurtownie danych to pomysł nowy, warto poświęcić nieco uwagi analizie wyzwań natury technologicznej przed którymi stają przedsiębiorstwa chcące zaimplementować to rozwiązanie. Pierwsze wyzwanie jest związane z rozmiarem aktywnej hurtowni danych. Współczesne hurtownie danych z łatwością osiągają rozmiary rzędu terabajta danych. Zgodnie z najnowszym rankingiem opublikowanym przez Winter Corp. największa funkcjonująca hurtownia danych należy do Land Registry, angielskiej agencji rządowej zarządzającej nieruchomościami i gruntami w Anglii i Walii i posiada rozmiar 18.3 terabajta. Kolejne miejsca w rankingu zajmują BT Group (11.7 terabajta) i United Parcel Service (9.0 terabajta). Takie rozmiary świadczą nie tylko o ilości danych składowanych w hurtowni danych, ale również dają pojęcie o stopniu skomplikowania schematu takiej hurtowni. Nierzadko na schemat składają się setki i tysiące tabel, tysiące wierzchołków referencyjnych i indeksów, ogromna liczba perspektyw i migawek. Efektywne zarządzanie tak ogromną bazą danych stanowi wielkie wyzwanie i poważny problem. Drugi czynnik technologiczny związany jest z dziedzina zastosowań aktywnych hurtowni danych. Twórcy hurtowni muszą być przygotowani na to, że gwałtowny wzrost liczby użytkowników pociąga za sobą lawinowy wzrost liczby wymagań stawianych aktywnej hurtowni danych. Użytkownicy, w tym przypadku bardziej zróżnicowani niż użytkownicy tradycyjnej hurtowni danych, formułują różne, czasami sprzeczne, wymagania odnośnie funkcjonalności hurtowni danych. Zaspokojenie wszystkich oczekiwań użytkowników może nie okazać się zadaniem trywialnym. Wreszcie trzecie wyzwanie związane jest z przeniesieniem ciężaru wykorzystania aktywnej hurtowni danych z biur analityków (back-office) do biur pracowników bezpośrednio obsługujących klientów (front-office). Ta zmiana pociąga za sobą konieczność zapewnienia wsparcia informatycznego dużej liczbie użytkowników którzy są skupieni na bieżącej działalności przedsiębiorstwa. Wsparcie musi też być udzielane szybciej, ponieważ od szybkości i jakości działania użytkowników operacyjnych zależy zadowolenie klienta.

Przykłady aktywnej hurtowni danych

Aby lepiej zobrazować pojęcie aktywnej hurtowni danych i wyraźnie zaznaczyć, czym ta architektura różni się od tradycyjnie stosowanych rozwiązań, poniżej zamieszczono analizę dwóch przykładów wdrożenia aktywnej hurtowni danych, które zakończyły się pełnym sukcesem.

3M

3M to firma o ponad stuletniej historii. Dzisiaj warta 16 miliardów dolarów, spółka ma swoje przedstawicielstwa w ponad 60 krajach i obsługuje klientów z 200 krajów, sprzedając produkty z katalogu liczącego ponad pół miliona pozycji. W latach 90-tych 3M posiadała 40 oddzielnych departamentów i 60 narodowych oddziałów. Każda jednostka samodzielnie zbierała informacje o klientach, produktach, dostawcach, oraz sprzedaży. Istniały pewne aplikacje o charakterze analitycznym, ale pozwalały one na niewielką integrację informacji z kilku połączonych jednostek. Nie istniał żaden globalny widok danych przedsiębiorstwa, a każda jednostka obsługiwała swoich

klientów tak, jak gdyby była niezależna firma. Na to nakładają się niespójność i niekompletność danych, skutecznie uniemożliwiające utworzenie globalnej perspektywy przedsiębiorstwa.

W 1996 firma zainicjowała ambitny projekt, którego efektem było powstanie aktywnej hurtowni danych. Celem projektu było zbudowanie spójnego globalnego obrazu każdego produktu, klienta i partnera handlowego. Rewolucyjna zmiana dotknęła także strukturę firmy, na miejsce 10 departamentów odpowiadających technologiom oferowanym przez 3M powołano 6 departamentów odpowiadających docelowym grupom klientów. Cała informacja o każdym indywidualnym kliencie została scentralizowana i wprowadzona do *globalnej hurtowni danych przedsiębiorstwa* (ang. global enterprise data warehouse). Utworzona aktywna hurtownia danych została otwarta w 1997 roku i od tego czasu stała się podstawowym narzędziem działania korporacji. Z technologicznego punktu widzenia najtrudniejszym zadaniem podczas konstrukcji aktywnej hurtowni danych okazała się integracja ogromnych wolumenów danych z setek niezależnych systemów i dostarczenie wyników do geograficznie rozproszonej grupy użytkowników liczącej tysiące użytkowników.

Dzisiaj aktywna hurtownia danych umożliwia powstawanie nowych aplikacji, zarówno operacyjnych jak i analitycznych. Aplikacje zarządzające planowaniem zapotrzebowania i zamówień, aplikacje logistyczne, magazynowe, systemy sprzedaży, narzędzia analizujące rentowność, narzędzia do przewidywania trendów, udostępnianie zasobów klientom i kontrahentom, każdy rodzaj aplikacji może powstać na bazie aktywnej hurtowni danych, ponieważ przedstawia ona sama jedyną i aktualną wersję prawdy. Utworzone aplikacje znacznie zwiększyły produktywność i efektywność procesów biznesowych, uproszczyły współpracę z kontrahentami, umożliwiły stworzenie zaawansowanego interfejsu webowego dla klientów. Pośród bezpośrednich korzyści przedstawiciele 3M wymieniają zwiększenie sprzedaży wiązanych oraz budowanie aktywnych związków z klientami, których zaufanie do firmy znacząco wzrosło. Sama tylko konsolidacja rozproszonych działów IT zaowocowała zmniejszeniem nadmiarowych kosztów, zlikwidowaniem konkurencyjnych projektów, oraz przyspieszeniem działania. Ocenia się, że w przypadku działów IT całej korporacji oszczędności wyniosły 100 milionów dolarów. Imponująco wygląda też liczba opisująca oszczędności 3M wynikające z automatyzacji logistyki, szacuje się ją na miliard dolarów w przeciągu całego życia projektu. Wszystkie operacje sprzedaży są przetwarzane przez jeden centralny system, który zapewnia każdemu z krajowych oddziałów firmy własny, indywidualny widok danych, troszcząc się o konieczne transformacje (np. przeliczanie walut). Nawet takie zadania jak drukowanie etykiet na paczki wysyłane do klientów poczta są obsługiwane przez aktywną hurtownię danych. Dane przechowywane w aktywnej hurtowni danych są na najwyższym poziomie szczegółowości. Odświeżanie danych odbywa się w niektórych podsystemach w sposób ciągły (tak jest w przypadku obsługi zamówień), lecz niektóre podsystemy stosują odświeżanie wsadowe z częstotliwością co dwie godziny (zamówienia kierowane do fabryk).

Ford

Ford Motor Company to druga co do wielkości samochodowa firma świata. Liczy 350 000 pracowników a sprzedaż sięga 160 milionów dolarów rocznie. Bezpośrednio z firmą kontaktuje się 2000 dostawców, 6000 przedstawicieli, firma zarządza też 8 centrami dystrybucji w których pracuje 4000 pracowników. Centra dystrybucji zajmują się zaopatrzeniem 50 milionów samochodów w 200 000 różnych rodzajów części. Zarządzaniem dostawami części zajmuje się aplikacja IMAS (Inventory Management and Alerting System). Aplikacja ta była początkowo pomyślana jako proste narzędzie raportujące, którego zadaniem było informowanie użytkowników o ciężarówkach przewożących części, które z jakiegoś powodu były opóźnione. Ewolucja wymagań sprawiła, że aplikacja ta zaczęła coraz bardziej zmierzać w kierunku aktywnej hurtowni danych. Dzisiaj IMAS to operacyjna składnica danych posiadająca rozległe możliwości operacyjne i analityczne. System co noc mierzy poziom zaopatrzenia każdego centrum dystrybucji w części i śledzi geograficzne położenie każdego samochodu dostawczego. Analiza tych danych pozwala na codzienne monitorowanie i modyfikowanie przepływu części i dynamiczną zmianę priorytetów i punktów docelowych poszczególnych samochodów dostawczych. W ten sposób ciężarówka opuszczająca fabrykę w Detroit z najniższym priorytetem po dotarciu do Seattle może posiadać już pierwszeństwo rozładunku, jeżeli wiezione przez nią części są aktualnie potrzebne. IMAS pozwala nie tylko na monitorowanie samochodów dostawczych, ale również wspomaga pracę magazynierów w centrach zaopatrzenia. Każda część jest skanowana i system natychmiast podpowiada pracownikowi, gdzie dana część powinna się znaleźć. Należy zauważyć, że jest to ekstremalny przykład rozszerzenia grona użytkowników aktywnej hurtowni danych, ponieważ w tym przypadku użytkownikami są pracownicy fizyczni zajmujący się zaopatrzeniem magazynu w części. Z technologicznego punktu widzenia dwa wyzwania wydają się najbardziej interesujące: zapewnienie natychmiastowej odpowiedzi na zapytanie skierowane do aktywnej hurtowni danych (natychmiast po zeskanowaniu kodu kreskowego części system musi odpowiedzieć, w którym magazynie i na której półce położona jest dana część) oraz zmiana częstotliwości odświeżania hurtowni danych z cyklu miesięcznego na dzienny.

Integracja aplikacji przedsiębiorstwa

Zaawansowana architektura aktywnej hurtowni danych wymaga daleko idącej integracji aktualnie istniejących systemów analitycznych z systemami operacyjnymi. *Integracja aplikacji przedsiębiorstwa* (ang. Enterprise Application Integration, EAI) zapewnia poprawną i efektywną propagację danych między tymi systemami i dostarcza platformę, na której integracja heterogenicznych systemów może odbywać się możliwie najmniejszym kosztem. Integracja aplikacji przedsiębiorstwa oznacza zapewnienie nieograniczonej wymiany informacji pomiędzy wszystkimi komponentami składającymi się na infrastrukturę informatyczną przedsiębiorstwa [5]. Taka integracja jest bardzo trudnym zadaniem, wymaga bowiem stworzenia infrastruktury technicznej, która umożliwi płynne połączenie procesów biznesowych, platform programowych i sprzętowych, czy standardów i architektur informatycznych. W tym artykule koncentrujemy się na roli, jaką gra integracja aplikacji przedsiębiorstwa w stworzeniu aktywnej hurtowni danych.

Aktywna hurtownia danych wymaga ciągłego odświeżania strumieniem nowych danych pochodzących z systemów operacyjnych. Proces odświeżania nie jest trywialny i wymaga efektywnych narzędzi informatycznych. Jednym z takich narzędzi jest właśnie integracja aplikacji przedsiębiorstwa. Poprawnie zaimplementowana, umożliwia na odświeżanie zawartości aktywnej hurtowni danych w czasie rzeczywistym lub zbliżonym do rzeczywistego. Z drugiej strony, integracja aplikacji przedsiębiorstwa umożliwia także połączenie w drugą stronę, tzn. dostarczanie wyników analiz przeprowadzonych w komponencie analitycznym do aplikacji operacyjnych. Taka infrastruktura zapewnia sprzężenie zwrotne w procesie podejmowania decyzji: decyzje są podejmowane na podstawie integralnych i aktualnych danych, podjęta decyzja jest natychmiast transmitowana do aplikacji operacyjnej i na jej podstawie podejmowana jest bezzwłocznie akcja. Komunikacja między źródłami danych operacyjnych, aplikacjami analitycznymi i aplikacjami operacyjnymi odbywa się za pomocą specjalnej *magistrali komunikatów* (ang. EAI message bus).

Skuteczna implementacja integracji aplikacji przedsiębiorstwa wymaga spełnienia kilku warunków. Dane transmitowane magistralą pomiędzy poszczególnymi komponentami mają często charakter poufny, system przesyłania wiadomości musi więc charakteryzować się bezpieczeństwem, niezawodnością, wydajnością i skalowalnością. Komunikujące się ze sobą moduły muszą być oczywiście połączone siecią komputerową. Ponieważ przedsiębiorstwa implementujące takie rozwiązanie to najczęściej duże, rozproszone geograficznie organizmy, integracja aplikacji przedsiębiorstwa musi brać pod uwagę możliwości i zagrożenia związane z wykorzystaniem sieci rozległych do komunikacji między poszczególnymi komponentami architektury. Magistrala przesyłająca komunikaty musi posiadać zdolność do zarządzania przepływem komunikatów, co pociąga za sobą umiejętność konwersji komunikatów. Integrowane moduły mogą być wysoce heterogeniczne, a koszty ich modyfikacji wysokie, zatem sposób włączania danego modułu do magistrali wymiany komunikatów musi wiązać się z możliwie najmniejszymi kosztami. W idealnym przypadku włączenie danego modułu do magistrali wymiany komunikatów nie powinno się wiązać z żadną ingerencją w kod aplikacji.

Technologie umożliwiające zbudowanie architektury integrującej aplikacje przedsiębiorstwa można podzielić na następujące klasy [2]. W ogólności każda kolejna klasa stanowi bardziej zaawansowaną implementację oferującą zwiększone możliwości.

- Integracja na poziomie danych (data-level integration): w tym przypadku integracja między aplikacjami osiągnięta jest poprzez współdzielenie danych. Dane są konwertowane na format odpowiedni dla każdej aplikacji i ładowane za pomocą tradycyjnych technik ETL i ELT. Taka formuła integracji nie pociąga za sobą praktycznie żadnych kosztów związanych z przystosowaniem istniejących aplikacji. Istotne jest, aby system zapewniał okna bezczynności umożliwiające przesłanie danych między poszczególnymi komponentami. Wada rozwiązania jest to, że dane w żadnym z integrowanych systemów nigdy nie są w pełni aktualne.
- Integracja na poziomie komunikatów (message-level integration): takie rozwiązanie umożliwia współdzielenie danych w czasie rzeczywistym. Aplikacje komunikują się ze sobą wykorzystując oprogramowanie do przesyłania wiadomości (IBM MQSeries, Microsoft MSMQ), które umożliwia zarówno komunikację typu point-to-point, jak i subskrypcję kolejek wiadomości. Wadą tego rozwiązania jest to, że komunikujące się komponenty muszą być świadome istnienia magistrali wiadomości, tzn. każdy komponent musi implementować interfejsy do wysyłania i odbierania komunikatów. To z kolei może pociągać za sobą konieczność ingerencji w kod istniejących już aplikacji.
- Integracja na poziomie procesów (process-level integration): jest to rozszerzenie wcześniejszego rozwiązania o dodanie zdolności do *zarządzania przepływem pracy* (ang. workflow management) i przepływem procesów między komunikującymi się komponentami. Rozwiązanie to, mimo że wymagające największych inwestycji, zapewnia też najbardziej zaawansowaną i daleko idącą integrację aplikacji przedsiębiorstwa.

Stworzenie infrastruktury umożliwiającej integrację aplikacji przedsiębiorstwa jest pierwszym krokiem do stworzenia aktywnej hurtowni danych. Oczywiście, sama tylko możliwość współdzielenia danych między systemami działającymi w przedsiębiorstwie to jeszcze nie wszystko. Konieczne jest zapewnienie sprzężenia zwrotnego między aplikacjami analitycznymi i operacyjnymi, wzbogacenie aplikacji operacyjnych o możliwość

korzystania z decyzji generowanych w aplikacjach analitycznych, i wiele innych. Niemniej jednak, integracja aplikacji przedsiębiorstwa stanowi warunek konieczny stworzenia aktywnej hurtowni danych.

Podsumowanie

W niniejszym artykule przedstawiono pojęcie operacyjnej składnicy danych. Wskazano cechy odróżniające ją od zarówno tradycyjnej hurtowni danych, jak i operacyjnej bazy danych. Omówiono zostały metody odświeżania operacyjnej składnicy danych. Operacyjna składnica danych jest podstawowym komponentem aktywnej hurtowni danych. Pojęcie aktywnej hurtowni danych zostało przedstawione szczegółowo, wraz z prezentacją korzyści wynikających z zaimplementowania tego narzędzia informatycznego. Kolejno przedstawiono czynniki wpływające na udaną implementację oraz podkreślono wyzwania technologiczne wynikające z wdrożenia aktywnej hurtowni danych. Dwa przykłady wdrożeń aktywnej hurtowni danych miały na celu przekonanie czytelnika o niebagatelnych zyskach płynących z implementacji aktywnej hurtowni danych. Artykuł kończy dyskusja nad metodami umożliwienia integracji aplikacji przedsiębiorstwa, która to integracja jest niezbędna do stworzenia środowiska pozwalającego na wdrożenie aktywnej hurtowni danych.

Literatura

- [1] H. Garcia-Molina, J.D. Ullman, J. Widom, Database System Implementation, Prentice Hall, 2000
- [2] B. Gold-Bernstein, EAI Market Segmentation, EAI Journal, Jul/Aug 1999
- [3] R. Hackathorn, Value Proposition for the Active Data Warehousing: Achieving the Intelligent Enterprise, Feb 2002
- [4] M. Jarke, M. Lenzerini, Y. Vassiliou, P. Vassiliadis, Fundamentals of Data Warehouses, Springer-Verlag Berlin Heidelberg, 2000
- [5] D. Linthicum, Enterprise Application Integration, Addison-Wesley, 1999