

## Systemy rozproszonych baz danych – 2

Fragmentacja, replikacja, zarządzanie transakcjami

Wykład przygotował:

Robert Wrembel



## Plan wykładu

- Cel i techniki fragmentacji danych
- Cel i techniki replikacji danych
- Problematyka zarządzania transakcjami rozproszonymi

ZSBD – wykład 2 (2)

Celem wykładu jest omówienie podstawowych zagadnień związanych z implementacją systemu rozproszonej bazy danych. Omówione zostaną następujące zagadnienia:

- cel i techniki fragmentacji danych (tzw. partycjonowanie),
- cel i techniki replikacji danych,
- problematyka zarządzania transakcjami rozproszonymi.

Problematyka optymalizacji zapytań rozproszonych zostanie omówiona w ramach laboratorium nr 3, wraz z praktyczną ilustracją konkretnych technik w systemie Oracle9i/10g.



## Wprowadzenie

- Problematyka projektowania RBD, optymalizacja zapytań, zarządzanie współbieżnością transakcji ⇒ znacznie trudniejsze niż w klasycznych BD
- Systemy komercyjne
  - Oracle9i/10g
  - Sybase Adaptive Server Enterprise, Adaptive Server Anywhere,
  - IBM DB2
  - MS SQLServer2000, SQLServer2005

ZSBD – wykład 2 (3)

Wiele problemów związanych z projektowaniem i zarządzaniem scentralizowanymi bazami danych, m.in. projektowanie struktury bazy danych, przetwarzanie i optymalizacja zapytań, zarządzanie współbieżnością transakcji, staje się znacznie trudniejsze w przypadku baz rozproszonych.

Dostępne na rynku komercyjne SZBD (Oracle9i/10g, Sybase Adaptive Server Enterprise i Sybase Adaptive Server Anywhere, IBM DB2, MS SQLServer2000 i SQLServer2005) oferują mniej lub bardziej zaawansowaną funkcjonalność wymaganą przy budowie systemów RBD.



## Fragmentacja/partycjonowanie (1)

- Podział tabeli na części, zwane fragmentami lub partycjami
- Techniki fragmentacji
  - pozioma
  - pionowa
  - mieszana

ZSBD – wykład 2 (4)

Jednym z podstawowych zagadnień projektowania RBD jest określenie sposobu podziału danych pomiędzy węzły, tak aby efektywność całego systemu była zadowalająca. Jedną z technik stosowanych w tym zakresie jest fragmentacja i alokacja.

Fragmentacja polega na podziale obiektu przechowującego dane, najczęściej tabeli, na mniejsze części, zwane fragmentami lub partycjami. W praktyce wyróżnia się trzy techniki fragmentacji, tj. fragmentację poziomą, pionową i mieszaną.



## Fragmentacja/partycjonowanie (2)

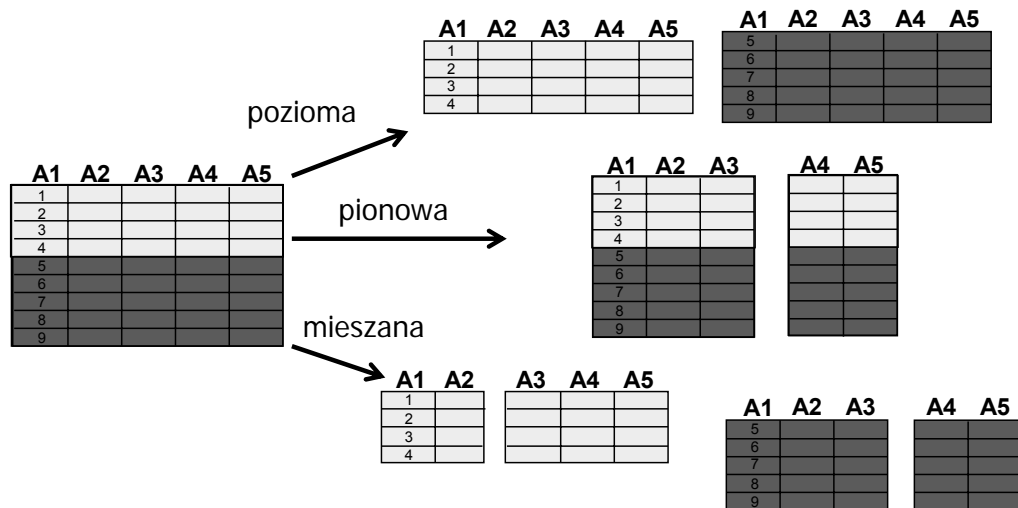
- Cel: zwiększenie efektywności dostępu do danych
  - zmniejszenie rozmiaru danych, które należy przeszukać
  - zrównoleglenie operacji dostępu do dysków, na których umieszczono fragmenty
  - alokowanie fragmentów "blisko" miejsca ich wykorzystania - redukcja kosztów transmisji sieciowej

ZSBD – wykład 2 (5)

Celem fragmentacji jest zwiększenie efektywności dostępu do danych. Dzięki właściwemu doborowi kryterium podziału obiektów na fragmenty można zapewnić, że aplikacje użytkowników będą adresowały tylko wybrane interesujące je partycje, w ten sposób ograniczając wolumeny przeszukiwanych danych. Ponadto, każdy z fragmentów może być umieszczony na innym dysku, co zapewni równoległy dostęp do wielu fragmentów równocześnie. Jeżeli fragmenty zostaną umieszczone "blisko" miejsca ich wykorzystania, wówczas można zredukować koszty transmisji sieciowej w sieciach o niskiej przepustowości.



## Rodzaje fragmentacji



ZSBD – wykład 2 (6)

Na slajdzie w sposób schematyczny przedstawiono wspomniane wcześniej techniki fragmentacji.



## Fragmentacja pozioma

- Podział zbioru rekordów na podzbiory
- Każdy podzbiór opisany identyczną liczbą atrybutów
- Wybór fragmentu (partycji) do którego trafi rekord realizowany na podstawie wartości jednego lub kilku wybranych atrybutów tabeli – tzw. atrybutów fragmentujących (partycjonujących)

ZSBD – wykład 2 (7)

Fragmentacja pozioma (ang. horizontal fragmentation) umożliwia podział zbioru rekordów tabeli na mniejsze podzbiory, z których każdy jest opisany identyczną liczbą atrybutów.

Wybór fragmentu (partycji) do którego trafi rekord jest realizowany na podstawie wartości jednego lub kilku wybranych atrybutów tabeli – tzw. atrybutów fragmentujących (partycjonujących).



## Fragmentacja pionowa

- Rozbicie tabeli na fragmenty złożone z podzbioru atrybutów
- Każdy fragment zawiera identyczną liczbę rekordów
- Atrybut niekluczowy X może się znaleźć tylko w jednym fragmencie
- Atrybut kluczowy znajduje się w każdym fragmencie  
– służy do łączenia fragmentów

ZSBD – wykład 2 (8)

Fragmentacja pionowa (ang. vertical fragmentation) umożliwia podział tabeli w pionie, tj. jej rozbicie na fragmenty złożone z podzbiorów atrybutów tabeli pierwotnej. Każdy fragment zawiera identyczną liczbę rekordów. Dany atrybut może się znaleźć tylko w jednym fragmencie. Nie dotyczy to atrybutów kluczowych, które są umieszczane w każdym fragmencie i służą one do łączenia fragmentów.





## Fragmentacja mieszana

- Połączenie fragmentacji poziomej i pionowej
- Warianty
  - Fragmentacja pozioma + pionowa
  - Fragmentacja pionowa + pozioma

ZSBD – wykład 2 (9)

Fragmentacja mieszana (ang. hybrid fragmentation) stanowi połączenie fragmentacji poziomej i pionowej. Występuje w dwóch wariantach. W pierwszym tabela jest najpierw dzielona poziomo, a następnie wszystkie lub wybrane jej fragmenty są dalej dzielone pionowo. W wariacie drugim, tabela jest najpierw dzielona pionowo, a następnie wszystkie lub wybrane jej fragmenty są dalej dzielone poziomo.



## Kryteria poprawności fragmentacji (1)

- Kompletność (ang. *completeness*)
  - jeżeli tabela  $T$  została podzielona na partycje  $TP_1, TP_2, \dots, TP_n$ , to każdy rekord z  $T$  lub jego fragment musi się znaleźć w jednej z partycji  $TP_1, TP_2, \dots$ , lub  $TP_n$
  - żadne dane nie są tracone

ZSBD – wykład 2 (10)

Poprawna fragmentacja musi spełniać trzy następujące kryteria: kompletność, rozłączność i rekonstrukcję.

Kompletność (ang. *completeness*) oznacza, że jeżeli tabela  $T$  została podzielona na partycje  $TP_1, TP_2, \dots, TP_n$ , to każdy rekord z  $T$  lub jego fragment musi się znaleźć w jednej z partycji  $TP_1, TP_2, \dots$ , lub  $TP_n$ . Kryterium to gwarantuje, że na skutek partycjonowania żadne dane z tabeli pierwotnej  $T$  nie zostaną utracone.



## Kryteria poprawności fragmentacji (2)

- Rozłączność (ang. disjointness)
  - jeżeli tabela  $T$  została podzielona na partycje  $TP_1, TP_2, \dots, TP_n$ , to każdy rekord z  $T$  lub jego fragment nie może się znaleźć w więcej niż jednej partycji
  - nie pojawiają się dane nadmiarowe
  - wyjątek: partycjonowanie pionowe
    - klucze podstawowe występują w każdej partycji

ZSBD – wykład 2 (11)

Rozłączność (ang. *disjointness*) oznacza, że jeżeli tabela  $T$  została podzielona na partycje  $TP_1, TP_2, \dots, TP_n$ , to każdy rekord z  $T$  lub jego fragment nie może się znaleźć w więcej niż jednej partycji. Kryterium to gwarantuje, że na skutek partycjonowania w bazie danych nie pojawiają się dane nadmiarowe. Wyjątkiem od tej reguły jest partycjonowanie pionowe, w którym wartości atrybutów stanowiących klucz podstawowy tabeli występują w każdej partycji.



## Kryteria poprawności fragmentacji (3)

- Rekonstrukcja (ang. *reconstruction*)
  - możliwość zrekonstruowania pierwotnej tabeli  $T$  ze wszystkich jej partycji
  - nie prowadzi do utraty danych
  - nie prowadzi do powstania danych nadmiarowych
- Rekonstrukcja z fragmentów poziomych
  - operator sumy zbiorów
- Rekonstrukcja z fragmentów pionowych
  - operator połączenia

ZSBD – wykład 2 (12)

Rekonstrukcja (ang. *reconstruction*) oznacza, że musi istnieć możliwość zrekonstruowania pierwotnej tabeli  $T$  ze wszystkich jej partycji. Rekonstrukcja ta nie może doprowadzić ani do utraty żadnych danych, ani do powstania danych nadmiarowych.

W przypadku partycjonowania poziomego, rekonstrukcji pierwotnej tabeli dokonuje się z wykorzystaniem operatora sumy zbiorów rekordów znajdujących się we wszystkich partycjach. W przypadku partycjonowania pionowego, rekonstrukcję realizuje się za pomocą łączenia partycji wykorzystując do tego klucz podstawowy tabeli pierwotnej.



## Algorytmy fragmentacji poziomej

- Round-robin
  - cykliczne rozmieszczanie danych w węzłach systemu
- Bazujący na wartości
  - rozmieszczenie danych w węzłach zależy od wartości danych
    - zakresowy
    - haszowy

ZSBD – wykład 2 (13)

Fragmentację poziomą implementuje się z wykorzystaniem dwóch podstawowych algorytmów, tj. round-robin i bazującego na wartości. Algorytm round-robin rozmieszcza rekordy cyklicznie we wszystkich węzłach systemu. W przypadku algorytmów bazujących na wartości, tj. zakresowego i hash'owego, rozmieszczenie danych w sieci zależy od wartości samych danych.



## Algorytm round-robin

- Równomierne rozpraszanie danych w węzłach sieci
  - równoważenie obciążenia
- Dane rozpraszane w sposób przypadkowy
  - odnalezienie żądanych informacji wymaga przeszukania wszystkich węzłów

ZSBD – wykład 2 (14)

Algorytm round-robin umożliwia równomierne rozpraszanie danych w węzłach sieci. Przykładowo, jeśli w sieci znajdują się trzy węzły, to pierwszy rekord tabeli zostanie umieszczony w węźle pierwszym, drugi — w węźle drugim, trzeci rekord — w węźle trzecim, czwarty — znów w węźle pierwszym itp. Ponieważ dane są rozpraszane w sposób przypadkowy, więc odnalezienie żądanych rekordów wymaga przeszukania wszystkich węzłów, co jest wadą tego rozwiązania.



## Algorytm zakresowy

- Każdy węzeł składa się z danych o wartościach atrybutu fragmentującego z zadanego zakresu lub o zadanej wartości



ZSBD – wykład 2 (15)

W algorytmie zakresowym, każdy węzeł przechowuje dane o wartościach atrybutu fragmentującego z zadanego zakresu lub o zadanej wartości. Przykładowo, na slajdzie przedstawiono trzy węzły BD1, BD2, BD3. Pierwszy z nich zawiera fragment tabeli faktury, przechowujący faktury z lat od 2000 do 2002. Drugi węzeł zawiera fragment tabeli przechowujący faktury z lat 2003-2005, a trzeci - zawiera fragment tabeli przechowujący faktury z roku 2006.



## Algorytm haszowy

- Dane umieszczane w węzłach na podstawie wartości funkcji haszowej
  - argument wywołania  $\Rightarrow$  wartość atrybutu fragmentującego
  - wynik  $\Rightarrow$  adres węzła
- Umieszczanie w jednym węźle rekordów z różnych powiązanych tabel
  - efektywność operacji łączenia tabel

ZSBD – wykład 2 (16)

W przypadku algorytmu fragmentacji 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 umieszczony rekord. W celu odnalezienia żądanych rekordów SZBD wykorzystujemy tę samą funkcję haszową, która została wykorzystana do fragmentacji danych. Zaletą tej metody jest możliwość automatycznego umieszczania w tym samym węźle rekordów pochodzących z różnych, powiązanych z sobą tabel. W ten sposób zwiększa się efektywność wykonywania operacji łączenia tabel, gdyż łączone z sobą rekordy znajdują się w tym samym węźle.





## Problem alokacji (1)

- Określenie miejsca fizycznego składowania danych/fragmentów
  - umieszczenie danych we właściwych węzłach sieci - "blisko" miejsca ich wykorzystania
  - kryterium alokacji: maksymalizacja efektywności systemu RBD

ZSBD – wykład 2 (17)

Ponieważ dane w systemie RBD mogą być składowane na dowolnym jego węźle, więc zachodzi konieczność takiego rozmieszczenia danych (np. tabel, fragmentów, indeksów), który zapewni największą efektywność całego systemu. Jest to tzw. problem alokacji.



## Problem alokacji (2)

- Problem alokacji
  - dla danego zbioru fragmentów
$$F=\{F_1, F_2, \dots, F_n\}$$
i zbioru węzłów systemu
$$W=\{W_1, W_2, \dots, W_m\}$$
w których działają aplikacje
$$A=\{A_1, A_2, \dots, A_j\}$$
  - znaleźć optymalny przydział  $F_i$  ( $i=1..n$ ) do  $W_j$  ( $j=1..m$ )

ZSBD – wykład 2 (18)

Problem alokacji można zdefiniować następująco: dla danego zbioru fragmentów  $F=\{F_1, F_2, \dots, F_n\}$  i danego zbioru węzłów systemu RBD  $W=\{W_1, W_2, \dots, W_n\}$  w których działa zbiór aplikacji użytkowników  $A=\{A_1, A_2, \dots, A_j\}$  należy znaleźć optymalny przydział fragmentów do węzłów.



## Problem alokacji (3)

- Optimum jest definiowane w kontekście
  - łącznego kosztu
    - przechowywania  $F_i$  w węźle  $W_j$
    - odczytu  $F_i$  z węzła  $W_j$
    - zmodyfikowania  $F_i$  w  $W_j$
    - komunikacji z  $W_j$
  - efektywności systemu
    - mierzona jego przepustowością w każdym z węzłów

ZSBD – wykład 2 (19)

Optimum przydziału jest definiowane w kontekście łącznego kosztu i efektywności systemu. Łączny koszt uwzględnia:

- koszt przechowywania fragmentu  $F_i$  w węźle  $W_j$ ,
- koszt odczytu fragmentu  $F_i$  z węzła  $W_j$ ,
- koszt zmodyfikowania fragmentu  $F_i$  w węźle  $W_j$ ,
- koszt komunikacji z węzłem  $W_j$ .

Efektywność systemu jest natomiast mierzona przepustowością w każdym z jego węzłów.

W literaturze zaproponowano wiele algorytmów alokacji. Ich omówienie wykracza jednak poza zakres tego wykładu.



## Replikacja danych

- Tworzenie kopii danych pochodzących z jednego węzła i przechowywanie ich w innych węzłach
- Cele
  - przyspieszenie dostępu do danych wykorzystywanych często
    - słownik miast, województw, kodów adresowych
    - katalog produktów
  - równoważenie obciążenia węzłów
  - równoległe wykonywanie zapytań
  - zwiększenie efektywności dostępu do danych przez wybór repliki do której dostęp jest najszybszy

ZSBD – wykład 2 (20)

Drugim podstawowym zagadnieniem projektowania RBD jest określenie sposobu duplikowania danych w wielu węzłach systemu RBD. Duplikowanie takie będziemy dalej nazywali replikacją danych.

Replikacja polega na tworzeniu kopii danych pochodzących z jednego węzła i przechowywaniu tych kopii w innych węzłach.

Cele replikacji są następujące:

- przyspieszenie dostępu do danych poprzez ich umieszczenie w węzle, który z nich intensywnie korzysta; przykładami mogą być tabele słownikowe miast, województw i kodów adresowych w systemach ewidencji, katalogi produktów w systemach sprzedaży;
- równoważenie obciążenia węzłów poprzez powielanie intensywnie wykorzystywanych (odczytywanych) danych do innych węzłów;
- równoległe wykonywanie zapytań adresujących te same dane; w takim przypadku 2 zapytania odwołujące się do tej samej tabeli mogą zostać wykonane na dwóch różnych replikach tej tabeli;
- zwiększenie efektywności dostępu do danych przez wybór repliki, do której dostęp jest najszybszy.



## Problematyka replikacji

- Określenie jednostki replikacji
- Określenie ilości replikowanych danych
- Określenie momentu odświeżania
- Określenie sposobu odświeżania



## Określenie jednostki replikacji

- Wskazanie tabeli, pionowego lub poziomego fragmentu zbioru tabel, tabeli w momencie definiowania repliki
  - za pomocą zapytania SQL
  - za pomocą narzędzi dostarczanych przez producenta konkretnego systemu

ZSBD – wykład 2 (22)

Określenie jednostki replikacji polega na wskazaniu jak duża jest zawartość pojedynczej repliki. Replika może być zbudowana na pojedynczej tabeli, pionowym lub poziomym jej fragmencie, bądź na zbiorze tabel. Jednostkę replikacji określa się w momencie definiowania repliki albo za pomocą zapytania SQL albo za pomocą narzędzi programistycznych dostarczanych przez producenta konkretnego systemu.



## Określenie ilości replikowanych danych (1)

- Bazy w pełni replikowane
  - każdy węzeł zawiera wszystkie dane z pozostałych węzłów
  - cechy
    - wysoka efektywność wykonywania zapytań adresujących dane z innych węzłów
    - duży narzut czasowy na odświeżanie replik

ZSBD – wykład 2 (23)

W zakresie określenia ilości replikowanych danych wyróżnia się bazy danych w pełni replikowane (ang. fully replicated databases) i bazy danych częściowo replikowane (ang. partially replicated databases).

W bazach w pełni replikowanych każdy węzeł zawiera wszystkie dane z pozostałych węzłów. Bazy takie cechują się wysoką efektywnością wykonywania zapytań, ponieważ każde zapytanie może zostać wykonane lokalnie. Wadą tego typu baz danych jest duży narzut czasowy na odświeżanie replik, ponieważ jest ich zbyt dużo. W praktyce tego typu bazy danych stosuje się niezwykle rzadko.



## Określenie ilości replikowanych danych (2)

- Bazy częściowo replikowane
  - węzeł może zawierać wybrane dane z innych węzłów
  - cechy
    - mniejsza efektywność wykonywania zapytań adresujących dane z innych węzłów
    - mniejszy narzut czasowy na odświeżanie replik

ZSBD – wykład 2 (24)

W bazach danych częściowo replikowanych każdy węzeł może zawierać wybrane dane z innych węzłów. Tego typu bazy danych cechuje mniejsza efektywność od poprzednich ponieważ część zapytań będzie wymagała sięgnięcia do zdalnych węzłów. Mniejszy jest też narzut czasowy na odświeżanie replik, ponieważ replik tych jest mniej niż w bazie danych w pełni replikowanej.





## Określenie momentu odświeżania

- Synchronicznie
  - natychmiast po zmianie danych w źródle
  - po zatwierdzeniu transakcji w źródle
- Asynchronicznie
  - na żądanie
  - automatycznie z zadaniem okresem

ZSBD – wykład 2 (25)

Ważną cechą odświeżania repliki jest wybór momentu odświeżania. W praktyce stosuje się dwie techniki, tj. odświeżanie synchroniczne i odświeżanie asynchroniczne.

Pierwsza technika polega na propagowaniu danych ze źródła do repliki albo natychmiast po zmianie zawartości źródła albo po zatwierdzeniu transakcji w źródle. Drugie rozwiązanie jest stosowane w praktyce.

Odświeżanie asynchroniczne polega na propagowaniu danych do repliki po jakimś czasie od wprowadzenia zmian w źródle. Możliwe są tu dwa rozwiązania, albo odświeżanie na żądanie repliki, albo odświeżanie automatyczne z zadaniem okresem. W praktyce stosuje się oba rozwiązania.



## Określenie sposobu odświeżania

- Pełne
  - ze źródła do repliki przesyłane wszystkie dane
  - cechy
    - łatwe implementacyjnie
    - kosztowne
- Przyrostowe
  - ze źródła do repliki przesyłane dane zmienione
  - cechy
    - trudniejsze implementacyjnie
    - mniej kosztowne

ZSBD – wykład 2 (26)

Kolejną ważną cechą odświeżania repliki jest określenie sposobu jej odświeżania. W praktyce stosuje się dwie techniki, tj. odświeżanie pełne i odświeżanie przyrostowe.

Odświeżanie pełne polega na każdorazowym przesyłaniu ze źródła do repliki wszystkich danych, które replika udostępnia. Ta technika odświeżania charakteryzuje się łatwością implementacji i dużym kosztem odświeżania ze względu na dużą ilość przesyłanych danych.

Odświeżanie przyrostowe polega na przesyłaniu do repliki wyłącznie zmian dokonanych w źródle od czasu ostatniego odświeżenia. Ta technika cechuje się trudniejszą implementacją i niższymi kosztami odświeżania.



## Implementacja replikacji

- Replika - migawka - perspektywa zmaterializowana
  - obiekt bazy danych, którego definicja zawiera
    - moment odświeżania
    - sposób odświeżania
    - zakres replikowanych danych z tabel źródłowych (najczęściej specyfikacja za pomocą zapytania SQL)
  - systemowy proces odświeżający

ZSBD – wykład 2 (27)

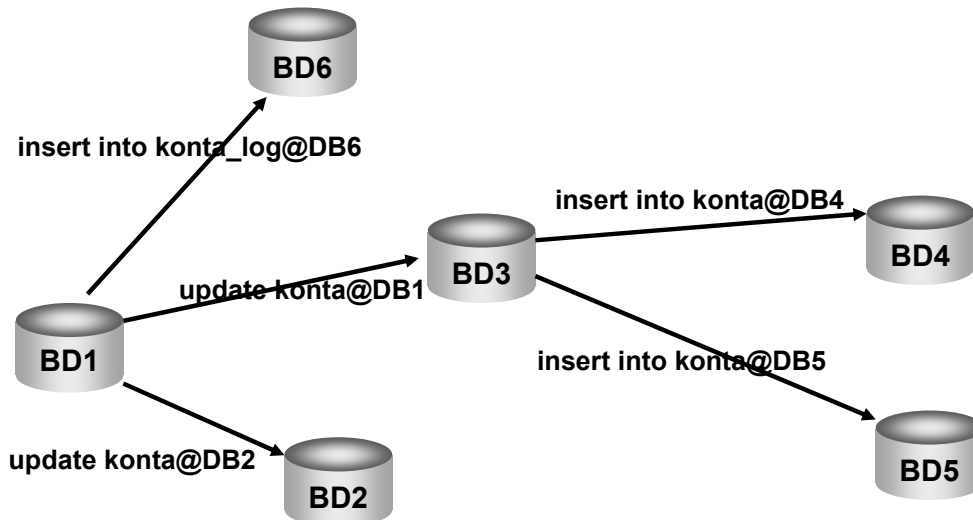
W praktyce replika jest implementowana jako pewien obiekt bazy danych, którego definicja zawiera:

- moment odświeżania,
- sposób odświeżania,
- zakres replikowanych danych z tabel źródłowych; zakres ten najczęściej specyfikuje się za pomocą zapytania SQL.

Z repliką jest związany systemowy proces odpowiedzialny za jej odświeżanie.



## Zarządzanie transakcjami w RBD



ZSBD – wykład 2 (28)

W systemie RBD użytkownik może wykonać transakcję, która adresuje wiele węzłów. Przykład takiej transakcji przedstawiono na slajdzie. W skład systemu wchodzi 6 baz danych BD1-BD6. Z bazy BD1 użytkownik wykonuje transakcję która wstawia dane do tabeli *konta\_log* w bazie BD6, modyfikuje zawartość tabeli *konta* w bazie BD2 i BD1. Ta ostatnia modyfikacja uruchamia polecenia *insert* do tabel w bazach DB4 i BD5. Transakcję taką będziemy nazywali rozproszoną (ang. distributed transaction) lub globalną (ang. global transaction).



## Transakcja rozproszona (1)

- Odwołuje się przynajmniej do jednego zdalnego węzła systemu RBD
- Transakcja rozproszona jest reprezentowana przez zbiór transakcji lokalnych
  - w każdej z baz danych, do której odwołuje się transakcja rozproszona tworzona jest jedna transakcja lokalna
- Cechy transakcji rozproszonej
  - trwałość
  - spójność
  - izolacja
  - atomowość

ZSBD – wykład 2 (29)

Transakcja rozproszona odwołuje się przynajmniej do jednego zdalnego węzła systemu RBD. Transakcja rozproszona jest reprezentowana przez zbiór **transakcji lokalnych**. W każdej z baz danych, do której odwołuje się transakcja rozproszona tworzona jest jedna transakcja lokalna. Zarówno każda z transakcji lokalnych jak i rozproszona mają cechy trwałości, spójności, izolacji i atomowości.



## Transakcja rozproszona (2)

- Atomowość
  - wszystkie transakcje lokalne zatwierdzone lub wszystkie wycofane
- Problemy
  - węzeł BD3 ulega awarii w trakcie realizowania polecenia INSERT
  - połączenie sieciowe z węzłem BD6 zostaje przerwane po wykonaniu polecenia INSERT
  - węzeł BD2 ulega awarii w trakcie zatwierdzania transakcji rozproszonej

ZSBD – wykład 2 (30)

Cecha atomowości w odniesieniu do transakcji rozproszonej oznacza, że wszystkie transakcje lokalne wchodzące w skład transakcji rozproszonej muszą zostać zatwierdzone. Jeśli choć jedna z transakcji lokalnych nie może zostać zatwierdzona, wówczas całą transakcję rozproszoną należy wycofać.

W czasie realizowania transakcji rozproszonej system RBD może ulec częściowej awarii. Przykładowo:

- węzeł BD3 (ze slajdu 28) może ulec awarii w trakcie realizowania polecenia INSERT,
- połączenie sieciowe z węzłem BD6 może zostać przerwane po wykonaniu polecenia INSERT,
- węzeł BD2 może ulec awarii w trakcie zatwierdzania transakcji rozproszonej.

W tych i wielu innych przypadkach system musi zagwarantować zakończenie rozproszonej transakcji, w taki sposób aby zagwarantować cztery jej omówione wcześniej cechy.



## Transakcja rozproszona (3)

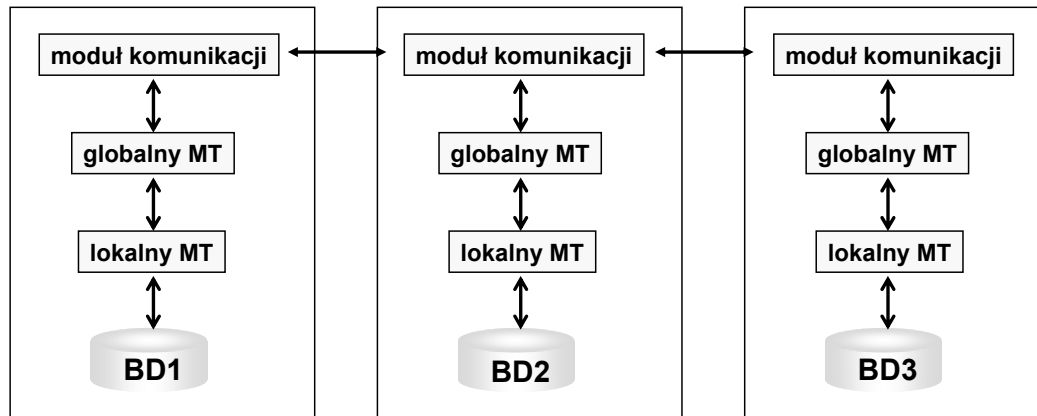
- Standardowy sposób zatwierdzania transakcji nie gwarantuje atomowości transakcji rozproszonej
- Potrzebny zaawansowany mechanizm zatwierdzania ⇨ protokół 2PC

ZSBD – wykład 2 (31)

Głównym problemem jest zagwarantowanie atomowości transakcji rozproszonej a standardowy mechanizm zatwierdzania transakcji nie gwarantuje jej atomowości. W związku z tym, zatwierdzanie lub wycofywanie transakcji rozproszonej, gwarantujące atomowość jest realizowane za pomocą specjalizowanego mechanizmu, tzw. **protokołu zatwierdzania dwu-fazowego** — *2PC* (ang. *two-phase commit*).



## Architektura zarządzania transakcjami rozproszonymi



ZSBD – wykład 2 (32)

Podstawową architekturę zarządzania transakcjami rozproszonymi przedstawiono na slajdzie. Każda z trzech baz danych BD1, BD2, BD3 posiada swój własny moduł *lokalnego menadżera transakcji* (lokalny MT), identycznie jak w standardowej scentralizowanej bazie danych. Ponadto, do zarządzania transakcjami rozproszonymi jest niezbędny moduł *globalnego menadżera transakcji* (globalny MT). Jego zadaniem jest koordynowanie wykonania zarówno lokalnych jak i rozproszonych transakcji zainicjowanych w jego węźle. Poszczególne węzły realizujące transakcję rozproszoną komunikują się za pośrednictwem *modułu komunikacji*, istniejącego w każdym węźle.





## Protokół 2PC (1)

- Fazy realizacji
  - przygotowanie/głosowanie (ang. voting phase)
  - decyzja (ang. decision phase)
- Warianty
  - scentralizowany 2PC
  - zdecentralizowany 2PC
  - liniowy 2PC

ZSBD – wykład 2 (33)

Protokół zatwierdzania 2-fazowego składa się z 2 faz: przygotowania, zwanej również głosowaniem (ang. voting phase) i decyzji (ang. decision phase).

Protokół ten może być implementowany w jednym z trzech wariantów:

- scentralizowanego 2PC,
- zdecentralizowanego 2PC,
- liniowego 2PC.



## Protokół 2PC (2)

- Węzły
- Koordynator globalny (KG)
  - węzeł inicjujący transakcję rozproszoną
  - koordynuje proces jej zatwierdzania/wycofywania
- Uczestnik (U)
  - węzeł realizujący transakcję lokalną będącą częścią rozproszonej

ZSBD – wykład 2 (34)

Z punktu widzenia protokołu 2PC w transakcji rozproszonej biorą udział dwa rodzaje węzłów, tj. koordynator globalny i uczestnik. Koordynator globalny jest węzłem inicjującym transakcję rozproszoną i koordynującym jej zatwierdzanie lub wycofywanie. Uczestnik jest węzłem realizującym transakcję lokalną, która to transakcja jest częścią transakcji rozproszonej.



## Scentralizowany 2PC (1)

- Koncepcja
  - KG pyta uczestników, czy są gotowi do zatwierdzenia
  - jeśli przynajmniej jeden uczestnik odpowiada ABORT, lub nie odpowie w zadanym czasie ⇒ KG wysyła do wszystkich uczestników komunikat wycofania transakcji
  - jeśli wszyscy uczestnicy gotowi ⇒ KG wysyła do nich komunikat zatwierdzenia transakcji

ZSBD – wykład 2 (35)

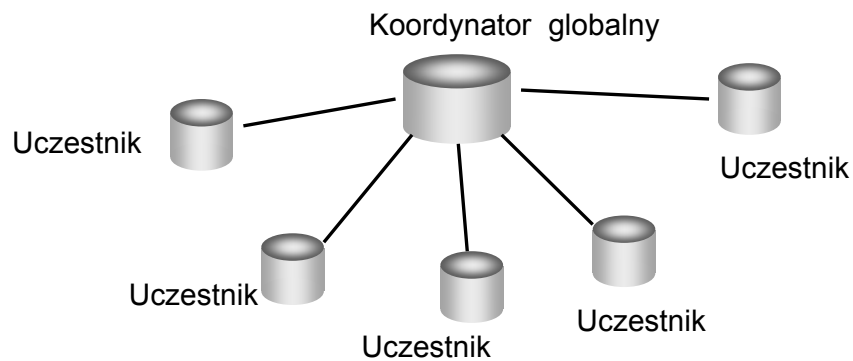
Koncepcja pracy protokołu 2PC jest następująca.

Koordinator globalny, pyta uczestników czy są gotowi do zatwierdzenia swoich transakcji lokalnych. Jeśli przynajmniej jeden uczestnik odpowiada komunikatem ABORT (nie jest gotowy), lub nie odpowie w zadanym czasie, wówczas koordinator globalny wysyła do wszystkich uczestników komunikat wycofania transakcji. Jeśli natomiast wszyscy uczestnicy zgłoszą gotowość do zatwierdzenia swoich transakcji lokalnych, wówczas koordinator globalny wysyła do nich komunikat żądający zatwierdzenia transakcji.



## Scentralizowany 2PC (1)

- Topologia komunikacyjna



ZSBD – wykład 2 (36)

Jak wspomniano, protokół 2PC może być implementowany w jednym z trzech wariantów:

- scentralizowanego 2PC,
- zdecentralizowanego 2PC,
- liniowego 2PC.

Topologia komunikacyjna scentralizowanego protokołu 2PC została przedstawiona na slajdzie. Jak widać, wszystkie węzły komunikują się z centralnym węzłem, którym jest koordynator globalny.



## Scentralizowany 2PC - przygotowanie (1)

- Koordynator globalny
  - zapisanie w lokalnym logu (na dysku) informacji o rozpoczęciu zatwierdzania
  - wysłanie do uczestników komunikatu PREPARE (przygotowanie do zatwierdzania)
  - przejście do stanu oczekiwania na komunikaty od uczestników

ZSBD – wykład 2 (37)

Na kolejnych slajdach zostaną szczegółowo omówione poszczególne fazy protokołu 2PC w topologii scentralizowanej. W praktyce ta topologia jest wykorzystywana najczęściej.

W fazie przygotowania, koordynator globalny wykonuje następujące czynności:

- zapisuje na dysku w swoim lokalnym pliku logu rekord mówiący o rozpoczęciu zatwierdzania transakcji;
- następnie wysyła do uczestników żądanie przygotowania do zatwierdzania; jest to komunikat PREPARE;
- po wysłaniu tego komunikatu KG przechodzi w stan oczekiwania na komunikaty od uczestników.



## Scentralizowany 2PC - przygotowanie (2)

- Uczestnik
  - odebranie komunikatu PREPARE od KG
  - a) przygotowanie do zatwierdzenia
    - zapisanie do logu stanu transakcji i rekordu o gotowości do zatwierdzenia
    - wysłanie do KG komunikatu READY\_COMMIT
  - b) przygotowanie do wycofania
    - zapisanie do logu stanu transakcji i rekordu o wycofywaniu
    - wysłanie do KG komunikatu ABORT
  - przejście do stanu oczekiwania na kolejny komunikat od KG

ZSBD – wykład 2 (38)

W fazie przygotowania, uczestnik odbiera komunikat PREPARE od koordynatora globalnego i przygotowuje się do zatwierdzenia swojej transakcji lokalnej. Przygotowanie to polega na zapisaniu w swoim lokalnym pliku logu stanu transakcji i rekordu mówiącego o gotowości do zatwierdzenia. Następnie uczestnik wysyła do KG komunikat READY\_COMMIT, potwierdzający przygotowanie.

W przypadkach, gdy uczestnik nie może się przygotować do zatwierdzenia lub odbierze od KG komunikat ABORT (przygotowanie do wycofania), wówczas zapisuje w swoim logu stan transakcji i rekord mówiący o wycofywaniu transakcji. Następnie uczestnik wysyła komunikat ABORT do KG i przechodzi do stanu oczekiwania na kolejny komunikat.



## Scentralizowany 2PC - decyzja (1)

- Koordynator globalny
  - a) jeżeli wszyscy uczestnicy odpowiedzieli `READY_COMMIT`
    - zapis do logu informacji o zatwierdzeniu
    - wysłanie komunikatu `GLOBAL_COMMIT` do uczestników
  - b) jeżeli przynajmniej jeden U odpowiedział `ABORT`
    - zapis do logu informacji o wycofaniu
    - wysłanie `GLOBAL_ABORT` do uczestników
- oczekiwanie na potwierdzenia od uczestników

ZSBD – wykład 2 (39)

W fazie decyzji, koordynator globalny wykonuje następujące czynności.

Jeżeli wszystkie odebrane komunikaty to `READY_COMMIT`, wówczas KG zapisuje do swojego logu rekord mówiący o zatwierdzeniu transakcji rozproszonej. Następnie wysyła komunikat `GLOBAL_COMMIT` do uczestników.

Jeżeli przynajmniej jeden uczestnik odpowiedział komunikatem `ABORT`, wówczas KG zapisuje do swojego logu rekord mówiący o wycofaniu transakcji rozproszonej. Następnie wysyła komunikat `GLOBAL_ABORT` do uczestników.

Po wysłaniu komunikatu, KG przechodzi w stan oczekiwania na potwierdzenia uczestników.



## Scentralizowany 2PC - decyzja (2)

- Uczestnik
- jeśli odebrany GLOBAL\_COMMIT
  - zapis do logu informacji o zatwierdzeniu
  - zatwierdzenie transakcji lokalnej
  - zwolnienie blokad i zasobów systemowych
  - wysłanie potwierdzenia do KG
- jeśli odebrany GLOBAL\_ABORT
  - zapis do logu informacji o wycofaniu
  - wycofanie transakcji lokalnej
  - zwolnienie blokad i zasobów systemowych
  - wysłanie potwierdzenia do KG

ZSBD – wykład 2 (40)

W fazie decyzji uczestnik wykonuje następujące czynności.

Jeżeli odebrał komunikat GLOBAL\_COMMIT, wówczas:

- zapisuje w swoim logu rekord informujący o zatwierdzeniu swojej transakcji lokalnej,
- zatwierdza transakcję lokalną,
- zwalnia blokady założone przez transakcję lokalną,
- zwalnia zasoby systemowe przeznaczone do obsługi transakcji lokalnej,
- wysyła potwierdzenie zatwierdzenia transakcji do KG.

Jeżeli uczestnik odebrał komunikat GLOBAL\_ABORT, wówczas:

- zapisuje w swoim logu rekord informujący o wycofaniu swojej transakcji lokalnej,
- wycofuje transakcję lokalną,
- zwalnia blokady i zasoby systemowe,
- wysyła potwierdzenie wycofania transakcji do KG.





## Scentralizowany 2PC - decyzja (3)

- Koordynator globalny
  - odbiór wszystkich potwierdzeń od uczestników
  - zapis do logu informacji o zakończeniu transakcji

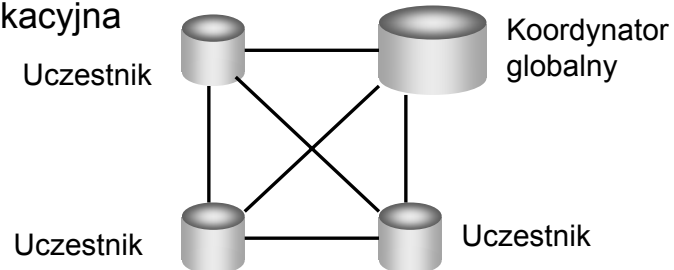
ZSBD – wykład 2 (41)

W ostatnim kroku fazy decyzji, KG odbiera potwierdzenia o zatwierdzeniu od wszystkich uczestników i zapisuje do logu rekord zakończenia transakcji.



## Zdecentralizowany 2PC (1)

- Topologia komunikacyjna



- Koncepcja

- KG wysyła komunikat `GLOBAL_COMMIT` (`GLOBAL_ABORT`) do wszystkich uczestników
- każdy uczestnik wysyła odpowiedź do wszystkich pozostałych węzłów

ZSBD – wykład 2 (42)

W architekturze zdecentralizowanej nie ma centralnego węzła - każdy węzeł komunikuje się z każdym.

Koncepcja działania protokołu 2PC w architekturze zdecentralizowanej jest następująca.

KG wysyła komunikat `GLOBAL_COMMIT` lub `GLOBAL_ABORT` do wszystkich uczestników. Uczestnicy przygotowują się w sposób identyczny do omówionego wcześniej, ale komunikat `READY_COMMIT` bądź `ABORT` wysyłają do wszystkich węzłów. W ten sposób, każdy węzeł ma pełen obraz stanu transakcji rozproszonej.



## Zdecentralizowany 2PC (2)

- Cechy
  - większy ruch sieciowy (większa liczba przesyłanych komunikatów)
  - każdy węzeł zna decyzje pozostałych węzłów i na tej podstawie zatwierdza/wycofuje sam
    - faza decyzji nie jest konieczna

ZSBD – wykład 2 (43)

Zdecentralizowany protokół 2PC cechuje się:

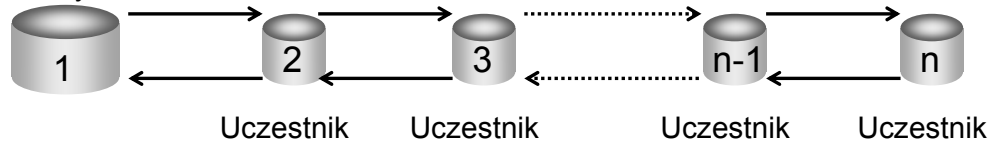
- większym ruchem sieciowym ze względu na większą liczbę przesyłanych komunikatów,
- możliwością wyeliminowania fazy decyzji ponieważ każdy uczestnik zna odpowiedź pozostałych uczestników i na jej podstawie może podjąć decyzję o zatwierdzeniu lub wycofaniu swojej transakcji lokalnej.



## Liniowy 2PC (1)

- Topologia komunikacyjna

Koordinator globalny



- Koncepcja
  - węzły są uporządkowane liniowo
  - każdy węzeł otrzymuje numer od 1 (KG) do n (U)
  - węzeł o numerze  $i$  komunikuje się tylko z węzłami sąsiednimi, tj. o numerach  $i-1$ ,  $i+1$

ZSBD – wykład 2 (44)

W architekturze liniowego 2PC węzły są uporządkowane liniowo i każdy z nich otrzymuje numer. Węzeł KG otrzymuje numer 1. Uczestnicy otrzymują kolejne numery, 2, 3, do n. W czasie wymiany komunikatów, każdy węzeł komunikuje się tylko z węzłami sąsiednimi, tj. węzłem bezpośrednio go poprzedzającym i węzłem bezpośrednio po nim następującym.



## Liniowy 2PC (2)

- Faza przygotowania
  - komunikaty wysyłane od węzła 1 do n
  - podejmuje decyzję, wysyła swoją decyzję z komunikatem do węzła o numerze  $i+1$
  - węzeł o numerze  $i+1$  podejmuje swoją decyzję i załącza ją w komunikacie do węzła  $i+2$ , itd.
  - ostatni węzeł w łańcuchu podejmuje decyzję o zatwierdzeniu/wycofaniu na podstawie zawartości komunikatu od węzła  $n-1$ 
    - komunikat ten zawiera decyzje wszystkich wcześniejszych węzłów

ZSBD – wykład 2 (45)

W fazie przygotowania, komunikaty są wysyłane przez KG do węzła o numerze 2. Węzeł ten po przygotowaniu się dokłada do komunikatu otrzymanego od KD swoją decyzję. Taki rozszerzony komunikat jest następnie wysyłany do węzła 3. Węzeł 3 po przygotowaniu się dokłada swoją decyzję do komunikatu otrzymanego od węzła 2 i wysyła tak rozszerzony komunikat do węzła 4. Każdy kolejny węzeł rozszerza otrzymany komunikat o swoją decyzję dopóki komunikat nie dotrze do węzła ostatniego. Po otrzymaniu komunikatu, węzeł ostatni podejmuje decyzję o zatwierdzeniu lub wycofaniu swojej transakcji. Decyzja ta jest podejmowana na podstawie zawartości komunikatu, który otrzymał. Zawiera on bowiem informacje o przygotowaniu lub nieprzygotowaniu się wszystkich wcześniejszych węzłów.



## Liniowy 2PC (3)

- Faza decyzji
  - komunikaty wysyłane od węzła  $n$  do 1
  - węzeł  $i$  podejmuje decyzję o zatwierdzeniu/wycofaniu i załącza ją w komunikacie do węzła  $i-1$
  - komunikat trafiający do KG zawiera decyzje wszystkich węzłów

ZSBD – wykład 2 (46)

Po podjęciu decyzji o zatwierdzeniu (wycofaniu) ostatni węzeł umieszcza tę decyzję w komunikacie wysyłanym do węzła poprzedniego (o numerze  $n-1$ ). W fazie decyzji komunikat wędruje od węzła ostatniego do KG. Każdy węzeł po drodze dopisuje do otrzymanego komunikatu swoją decyzję. W ten sposób, komunikat docierający do KG zawiera decyzje wszystkich węzłów. Na tej podstawie KG podejmuje decyzję o ostatecznym zatwierdzeniu lub wycofaniu transakcji. Decyzja ta jest wysyłana komunikatem do węzła nr 2 i propaguje się do węzła ostatniego.