



AUTOREFERAT ROZPRAWY DOKTORSKIEJ

mgr inż. Krzysztof Rączka

Technologia wieloagentowej integracji danych w rozproszonym systemie medycznym

Promotor:
Prof. dr hab. inż. Wiesław Wajs

Wstęp

W klasie systemów rozproszonych, charakteryzujących się wykorzystaniem dużej ilości danych oraz reżimem czasowym, występuje problem przepustowości sieci i problem określenia kolejności procesów związanych z pozyskaniem, replikacją, przesyłaniem i zagospodarowaniem danych przez poszczególne procesy. System medyczny, będący przedmiotem rozważań rozprawy, stanowi szczególny przypadek, w którym dane dynamiczne pochodzące z urządzeń pomiarowych oraz specjalistycznej aparatury medycznej są wykorzystywane przez moduły predykcyjne, dedykowane do wspierania decyzji lekarzy. Z domeny zastosowania wprost wynikają ścisłe uwarunkowania czasowe, wymagające zapewnienia efektywności rozwiązania.

Centralnym zagadnieniem dysertacji jest autorski podsystem replikacji, zrealizowany w technologii wieloagentowej, który wspiera efektywną integrację danych zarówno w aspekcie replikacji lokalnej, jak i sieciowej integracji danych. W pierwszym zagadnieniu rozważa się automatyzację doboru wydajnego zakresu pracy (parametrów buforujących oraz stopnia zrównoleglenia przetwarzania) oraz lokalizacji agentów (metody transportu danych), a także automatyzację monitorowania. Drugie z nich stanowi rozszerzenie i uogólnienie tematyki poszukiwania maksymalnego przepływu w sieci komputerowej o mechanizmy wspierające priorytety poszczególnych procesów.

Unikalne zapewnienie wsparcia dla priorytetów projektów zapewniające efektywność replikacji jest uważane przez autora za ważny element rozprawy. Rozwiązanie stanowi rozszerzenie algorytmów trasowania stosowanych w sieci Internet, gdzie następuje optymalizacja drogi przesyłu bez rozpatrywania priorytetów realizacji. Podejście proponowane w zagadnieniu replikacji lokalnej jest również nowatorskie.

Zaprojektowany oraz zrealizowany przez autora podsystem replikacji, wykorzystujący przedstawione pomysły został użyty do wykazania prawdziwości tezy niniejszej rozprawy, która brzmi następująco:

Możliwe jest zastosowanie technologii wieloagentowej w celu zapewnienia efektywności replikacji w sieciowym środowisku rozproszonym, charakteryzującym się uwarunkowaniami czasowymi oraz dużą ilością wykorzystywanych danych.

Z tezy wynikają następujące zadania badawcze i realizacyjne, wpływające bezpośrednio na efektywność replikacji danych:

- Replikacja lokalna
 - automatyczne wyznaczenie wartości parametrów zapewniających wydajną ekstrakcję oraz ładowanie danych,
 - dynamiczna alokacja zasobów poprzez automatyczny dobór ilości agentów dedykowanych do realizacji konkretnego etapu integracji danych,
 - wczesne ostrzeganie o sytuacjach wyjątkowych oraz detekcja sytuacji krytycznych, determinujących decyzję o celowości rozpoczęcia lub kontynuacji przetwarzania,
 - optymalizacja lokalizacji agentów wykonawczych oraz związanej z tym metody transportu danych.
- Sieciowa integracja danych
 - minimalizacja maksymalnego kosztu (czasu) przesyłu danych uwzględniając priorytety poszczególnych procesów,
 - minimalizacja całkowitego kosztu (czasu) przesyłu danych uwzględniając priorytety poszczególnych procesów.

Koncepcja systemu medycznego

Przedmiotem rozważań jest rozproszony system medyczny, dedykowany do wspierania decyzji lekarzy. System generuje sugestie oraz wskazówki odnośnie stosownych metod leczenia. W celu realizacji zadania niezbędne jest zbieranie, replikacja oraz predykcja danych, dotyczących opisu parametrów określających stan zdrowia pacjentów.

W ramach oddziałów szpitalnych funkcjonują dedykowane systemy medyczne, typowo posiadające własne bazy danych składujące informację lokalną. Bazy te zasilane mogą być poprzez dane wprowadzane ręcznie, jak również automatycznie z poszczególnych urządzeń medycznych. Niektóre z urządzeń umożliwiają wewnętrzne buforowanie danych, a także zdalny dostęp do nich. Aparatura taka może być traktowana jak niezależne podsystemy funkcjonalne.

Pozyskane dane przetwarzane są poprzez poszczególne moduły analityczne, a także predycyjne realizowane na zbiorze serwerów centralnych. Wyróżnia się dedykowane serwery bazodanowe, a także obliczeniowe. Zakłada się realizację każdej analizy przypadku chorobowego w postaci oddzielnego projektu. Strategia taka dopuszcza równoległe prowadzenie wielu niezależnych badań. W rezultacie konieczne jest tworzenie baz danych na żądanie. Każda z nich posiada prywatną kopię niezbędnych danych.

Przyjęta architektura wymusza replikację dużej ilości danych pomiędzy serwerami lokalnymi a centralną bazą danych oraz nakłada znaczące wymagania wydajnościowe na realizowany system. Konieczne jest wprowadzenie mechanizmów priorytetowego rozwiązywania potrzeb w ścisłym reżimie czasowym, a także automatyzacja zarządzania projektami.

Podstawowe komponenty logiczne rozproszonego systemu medycznego: moduły obliczeniowe, centralna baza danych, a także bazy lokalne wdrożone są na fizycznych platformach sprzętowych, takich jak: serwer, pojedynczy komputer lub mikroprocesor urządzenia pomiarowego. Projekty należą do poszczególnych użytkowników. Każdy projekt wykorzystuje pojedynczą instancję modułu obliczeniowego oraz prywatną kopię centralnej bazy danych.

Krytyczność czasowego pozyskania danych określa priorytet danego projektu, który znany jest w momencie początkowym. Priorytet ten zależy od czasu niezbędnego do przeprowadzenia obliczeń stosując zadany moduł obliczeniowy oraz skali czasu, do której odnoszą się wyniki obliczeń. Wygenerowanie wyników predykcji po czasie, którego ona dotyczy jest bezużyteczne, powoduje jedynie niepotrzebne obciążenie systemu.

Integracja danych w systemie odbywa się pomiędzy lokalnymi bazami danych, a bazą centralną poprzez interfejsy danych. Interfejs jest jednostką atomową określającą pożądany przepływ danych. Opisuje on sposób odczytu danych z bazy źródłowej oraz zapisu do tabeli docelowej centralnej bazy danych. Opcjonalnie może zawierać definicję transformacji, korekty, a także walidacji danych. Interfejs danych realizuje rolę filtra, dzięki któremu czytane są wyłącznie niezbędne dane z wymaganej grupy tabel źródłowych.

Realizacja interfejsów danych odbywa się poprzez sieć komputerową, której składnikami są fizyczne platformy sprzętowe. Faktyczna realizacja tych połączeń zależy od przynależności komponentów do określonych platform sprzętowych oraz połączeń pomiędzy nimi. Uzyskaną strukturę można opisać jako graf, którego węzłami są fizyczne platformy sprzętowe, natomiast krawędziami poszczególne segmenty sieci komputerowej.

Wydajność replikacji pomiędzy dwoma ustalonymi węzłami w różnych kierunkach może się różnić. Wartość ta zależy nie tylko od prędkości sieci komputerowej, lecz również od wydajności poszczególnych węzłów oraz stosowanego na nich oprogramowania. Fakt ten powoduje konieczność opisu zagadnienia z zastosowaniem grafu zorientowanego.

Założenia projektowe podsystemu replikacji

Głównym zadaniem podsystemu replikacji jest zapewnienie wydajnego dostarczania danych do modułów obliczeniowych. W celu efektywnej realizacji tego zadania niezbędna jest optymalizacja procesu integracji danych, którą można podzielić na następujące podproblemy:

- replikację lokalną,
- sieciową integrację danych.

Przez replikację lokalną rozumie się wykonanie interfejsu danych pomiędzy dwoma, bezpośrednio połączonymi węzłami sieci komputerowej. Problem sformalizować można następująco. Graf zorientowany określony jest poprzez parę uporządkowaną $G = (V, E)$, gdzie V jest skończonym zbiorem węzłów (wierzchołków), a E jest skończonym zbiorem łuków grafu. W przyjętym modelu realizacja każdego interfejsu danych odbywa się poprzez łuki grafu. Dana jest lokalna, źródłowa baza danych LBD_i w węźle V_i oraz baza docelowa LBD_j lub CBD_j w wierzchołku V_j . Istnieje dokładnie jeden (ustalony) łuk (i, j) w zbiorze E , przez który odbywa się replikacja danych. Na obu platformach sprzętowych dostępny jest również serwer FTP. Możliwa jest realizacja replikacji na komputerze źródłowym, docelowym, bądź obydwu naraz. Należy tak dobrać architekturę fizyczną interfejsu danych (miejsce replikacji oraz sposób transferu danych), parametry buforujące agentów wykonawczych oraz stopień zrównoleglenia przetwarzania, aby zmaksymalizować transfer danych (przepustowość) pomiędzy zadanymi komputerami.

Ogólniejszy problem sieciowej integracji danych, dzieli się na dwa zagadnienia:

- wyznaczenia maksymalnego przepływu,
- doboru optymalnej sekwencji decyzji wynikającej z realizacji sekwencji zadań.

Rozwiązanie kwestii replikacji lokalnej dostarcza niezbędne dane wejściowe do wyznaczenia maksymalnego przepływu w zadanej sieci komputerowej. Zakłada się, iż graf zbudowany jest z wierzchołków reprezentujących fizyczne platformy sprzętowe, a także krawędzi symbolizujących sektory replikacji lokalnej. W zagadnieniu przyjmuje się, iż znana jest (wyznaczona w poprzednim etapie) prędkość replikacji pomiędzy każdymi dwoma, połączonymi węzłami sieci. Celem tej fazy zagadnienia optymalizacyjnego jest dobór drogi przesyłu danych oraz określenie maksymalnego możliwego przepływu w sieci komputerowej.

Problem maksymalnego przepływu w sieci komputerowej podzielić można na zagadnienie przesyłu danych jednorodnych oraz danych o różnorodnej strukturze. Pierwsze z nich rozwiązać można poprzez różne warianty implementacyjne metody Forda Fulkersona. Najszybszy znany algorytm (push-relabel) zaproponowany został przez Goldberga i Tarjana. Dla danych o różnorodnej strukturze rozwiązanie wyznaczyć można stosując specjalizację metody sympleks lub znacznie szybsze algorytmy znajdujące rozwiązania przybliżone. Najlepszy do tej pory algorytm przedstawiony został przez Karakostasa. Cechuje się on brakiem zależności od ilości typów danych podczas znajdowania przepływu niejawnego oraz minimalną (wielomianowo-logarytmiczną) zależnością dla przepływu jawnego.

Wyznaczone rozwiązanie problemu sieciowej integracji danych nie jest wystarczające w rozpatrywanym systemie medycznym. Wynika to z założonej konieczności obsługi priorytetów realizacji poszczególnych zadań. Wynik drugiego etapu procesu, łącznie z założeniem znajomości zadanych priorytetów realizacji poszczególnych zadań (związanych z maksymalnym terminem dostarczenia danych) stanowi wejście do realizacji procesu kombinatorycznego, dedykowanego podjęciu optymalnej sekwencji decyzji dla zadanego wskaźnika jakości. Przyjmuje się dwa podstawowe kryteria uporządkowania: minimalizację maksymalnego kosztu, a także minimalizację całkowitego kosztu realizacji zadania.

Model formalny sieciowej integracji danych

Przyjmijmy model obliczania niemalejącego w czasie t kosztu operacji oznaczonej indeksem i w postaci

$$c_i(t) = a_i t + b_i$$

Operacja składa się z etapów niezbędnych do realizacji sieciowej integracji danych. W modelu tym b_i identyfikuje maksymalną przepustowość grafu. Wartość ta jest ograniczona do przedziału (0-1]. Parametr a_i także należy do przedziału (0-1] i określa priorytet realizacji i -tego procesu. Dla $i=1, 2, \dots, n$ problemem jest określenie kolejności realizacji procesów.

Zakłada się, iż dany jest system przyznawania priorytetów dla każdego projektu. Priorytet jest związany z granicznym czasem określonym dla projektu, który wynika z ostatecznego terminu podjęcia decyzji przez lekarza. Nieznaczną część czasu projektu przeznaczona jest na dostarczenie danych. Tak więc zadany priorytet projektu pośrednio wyznacza maksymalny dozwolony czas replikacji.

Wyznaczenie wartości obu parametrów następuje w zadanym przedziale czasu $[t_1, t_2)$. W określonym na nowo przedziale czasu, wartości parametrów funkcji kosztu a_i oraz b_i są wyznaczone ponownie. Oddzielnym problemem jest określenie przedziału czasu $[t_1, t_2)$. Przedział ten jest wyznaczany przez algorytm nadrzędny i wynika z terminowości zrealizowania zadań stawianych systemowi.

Istnieją dwa, odmienne sformułowania funkcji celu:

- minimalizacja maksymalnego kosztu

$$c_{\max} = \max_{i=1, \dots, n} (c_i(t))$$

- minimalizacja całkowitego kosztu

$$\sum c_i = \sum_{i=1}^n c_i(t)$$

W pierwszym przypadku celem jest takie określenie kolejności dostępu agentów do węzłów grafu, aby przy najmniejszym koszcie udostępnione zostały do replikacji zasoby skupione w węzłach. W drugim przypadku celem jest takie określenie kolejności realizacji zadań agentów, aby wszystkie zadania wykonane zostały w jak najmniejszym czasie. Pierwszy przypadek preferuje agenta o najwyższym priorytecie (funkcji celu). W drugim dąży się do realizacji wszystkich zadań w najmniejszym koszcie.

Optymalne uporządkowanie zadań zależy od trybu pracy systemu. Podczas stosowania strategii konkurencji (dla przyjętego modelu) dane jest ono zależnością: $a_1/b_1 \geq a_2/b_2 \geq \dots \geq a_n/b_n$. W przypadku kooperacji (sumy kosztów) możliwość wyznaczenia uporządkowania optymalnego zależy od relacji pomiędzy parametrami a i b modelu. W sytuacji, kiedy przepustowość sieci nie jest ograniczeniem systemu parametr b można pominąć, wówczas optymalne uporządkowanie wyraża się w postaci: $a_1 \leq a_2 \leq \dots \leq a_n$. Jeżeli wszystkie współczynniki a są jednakowe, optymalna sekwencja wynosi: $b_1 \leq b_2 \leq \dots \leq b_n$. W przypadku identyczności parametru b (konieczne jest wstępne uporządkowanie sekwencji według współczynnika a_i) pierwszym zadaniem w sekwencji optymalnej jest zadanie o współczynniku a_j , takim że:

$$t_0 \leq b \left(1 + \sum_{k=j+1}^n \prod_{l=j+1}^k (a_l + 1) \right) \quad \wedge \quad t_0 > b \left(1 + \sum_{k=j+2}^n \prod_{l=j+2}^k (a_l + 1) \right)$$

W najbardziej ogólnym przypadku, kiedy wartości parametrów a_i , b_i są różne (zakłada się wstępne uporządkowanie współczynnika a_i), jeśli dla pewnego i spełniona jest nierówność:

$$t_0(a_i + 1) + b_i + ((a_j + 1)b_i + b_j) \left(1 + \sum_{k=1}^{n-2} \prod_{l=1, l \neq i, l \neq j}^k (a_l + 1) \right) \\ < t_0(a_j + 1) + b_j + ((a_i + 1)b_j + b_i) \left(1 + \sum_{k=1}^{n-2} \prod_{l=1, l \neq i, l \neq j}^k (a_l + 1) \right)$$

to zadanie o współczynniku i przyjmuje się jako pierwsze w sekwencji suboptymalnej.

Realizacja wieloagentowej integracji danych

Rozproszony system medyczny cechują następujące atrybuty technologii wieloagentowej: ograniczone możliwości wykonawcze poszczególnych klas agentów, zdecentralizowane sterowanie oraz asynchroniczne przetwarzanie rozproszonych danych. Dodatkowo, zaimplementowane zostały wszystkie obowiązkowe cechy architektury wieloagentowej, takie jak: reaktywność, autonomia, zorientowanie na cel oraz czasowa ciągłość. Jednocześnie (w wyniku specyfiki zagadnień integracyjnych) aspekty komunikatywności oraz mobilności mają zasadnicze znaczenie głównie dla agentów zarządzających oraz replikujących.

Podstawowe typy agentów w systemie to: agent zarządzający, replikujący, grupa agentów wykonawczych (agent czytający, transformujący, korygujący, sprawdzający oraz zapisujący) oraz agent monitorujący (zmiany strukturalne, dostępność zasobów oraz poziom błędów).

Zasadniczym zadaniem agenta zarządzającego jest globalna optymalizacja ruchu sieciowego z uwzględnieniem priorytetów projektów. Agent przypisany jest do realizacji pojedynczego projektu. Reprezentuje on użytkownika wewnątrz systemu. Odpowiedzialny jest za całość operacji związanych z tworzeniem projektu oraz replikacją niezbędnych danych. Koordynuje realizację obliczeń, a także dostarcza wyniki do użytkownika systemu. Agent pracować może w dwóch trybach pracy: kooperacji lub konkurencji z innymi agentami zarządzającymi. Celem pierwszego z nich jest optymalizacja korzyści z systemu dla wszystkich jego użytkowników, drugiego, zapewnienie realizacji zadania pojedynczego użytkownika.

Agent replikujący obsługuje całość lokalnego procesu integracyjnego od ekstrakcji danych z systemu źródłowego, poprzez opcjonalne etapy transformacji, uzupełnienia i korekty informacji, następującej po nich walidacji danych, aż do kończącego przetwarzania ładowania do systemu docelowego. Podstawowym zadaniem agenta replikującego jest lokalna replikacja danych, czyli fizyczne kopiowanie pomiędzy zadanym zbiorem źródłowym, a docelowym poprzez ustalony segment sieci. Agent spełnia rolę kontenera wykonawczego dla dedykowanych agentów wykonawczych. Dodatkowo, w zakresie jego obowiązków znajduje się dobór wydajnej architektury przetwarzania, realizowany zgodnie z wybraną przez siebie strategią: po stronie systemu źródłowego, docelowego, lub w obu miejscach naraz.

Optymalizacja lokalnej replikacji danych realizowana jest przez grupę agentów wykonawczych oraz agenta replikującego. Agenci wykonawczy rezydują na fizycznej platformie sprzętowej dobranej przez agenta replikującego. Ich zadania stanowią elementarne etapy procesu integracyjnego, a możliwości decyzyjne ograniczone są do doboru:

- odpowiednich wartości parametrów buforujących,
- stopnia zrównoleglenia przetwarzania.

Realizacja bazuje na szablonie agenta wykonawczego określającego wzorzec zachowania wszystkich klas agentów roboczych. Optymalizacja odbywa się na podstawie aktywnej obserwacji bezpośredniego środowiska wykonawczego – elementu buforującego łączącego poszczególne etapy procesu. Stosowany algorytm jest wariantem przeszukiwania binarnego, zakres sterownia parametrami procesu oparty jest na wynikach doświadczalnych.

Podsumowanie

W rozprawie zaprezentowano rozwiązanie problemu zapewnienia efektywności replikacji danych w rozproszonym systemie medycznym. Przedstawiono zagadnienia zarówno lokalnej, jak i sieciowej integracji danych. W pierwszym z nich zastosowanie technologii wieloagentowej umożliwiło automatyczne sterowanie wielkością buforów, doбором wydajnego stopnia zrównoleglenia operacji, a także fizyczną lokalizacją przetwarzania. Drugie stanowi połączenie zagadnień doboru optymalnej drogi sieciowej z tematyką poszukiwania efektywnego uszeregowania projektów uwzględniając priorytety realizacji.

Badania przeprowadzone stosując autorski podsystem replikacji wykazały, iż czytanie danych stosując strategię buforowania może skutkować dziesięciokrotnie wyższą efektywnością, niż podczas rezygnacji z tego mechanizmu. W przypadku zapisu różnice są jeszcze większe i mogą osiągać wartość pięćdziesięciokrotną. Największy zmierzony przyrost prędkości związany z przetwarzaniem równoległym kształtuje się na poziomie 70%. Podczas stosowania strategii zrównoleglenia replikacji z jawną warstwą transportową teoretyczna możliwość przyrostu wydajności przewyższa 65%, jednocześnie zaobserwowano faktyczną poprawę efektywności w granicach 40%.

Autor wykazał, iż już podczas zastosowania prostych algorytmów samoczynnego doboru nastaw parametrów procesu uzyskiwana wydajność plasuje się na poziomie 90% najlepszych pomierzonych wartości. Wyniki te kilkakrotnie przewyższają prędkość uzyskaną podczas replikacji zrealizowanej przy zastosowaniu podejścia klasycznego, co obrazuje celowość wykorzystania technologii wieloagentowej w dziedzinie integracyjnej. Dodatkowo technologia ta umożliwiła naturalne wprowadzenie aktywnego monitorowania zmian systemów źródłowych oraz docelowych, badania dostępności zasobów systemowych, a także bieżącego statusu przetwarzania.

Ważna część pracy dotyczy sieciowej integracji danych. Autor wykazał możliwość wykorzystania technologii wieloagentowej do zbierania danych i zarządzania przesyłem danych w sieciowym środowisku rozproszonym w celu zapewnienia realizacji wybranych zadań w określonym terminie. Zaprezentowano metody optymalizacji uporządkowania projektów zarówno podczas stosowania strategii konkurencji – problem optymalizacji postaci $n \mid a_i t_{i-1} + b_i; a_i > 0; b_i > 0; t_0 \geq 0 \mid \min(\max(t_i))$, jak i kooperacji pomiędzy poszczególnymi agentami – problem $n \mid a_i t_{i-1} + b_i; a_i > 0; b_i > 0; t_0 \geq 0 \mid \min \sum t_i$.

W sytuacji konkurencji zaprezentowano wydajną metodę wyznaczania uporządkowania optymalnego. Podczas strategii kooperacji szybki dobór sekwencji optymalnej realizowany jest dla przypadków, kiedy czynnik odpowiedzialny za przepustowość sieci jest nieistotny, bądź identyczny dla wszystkich rozpatrywanych projektów, lub czynnik określający priorytet wykonania jest taki sam. Dla najbardziej ogólnego przypadku, kiedy czynniki przyjętego modelu są różne, przedstawiono metodą znajdowania sekwencji suboptymalnej.

Wszystkie wykorzystane techniki zapewniające wydajność replikacji danych stanowią własny dorobek autora. Proponowane metody optymalizacji nie są prezentowane w żadnym z licznych projektów integracyjnych zebranych na Uniwersytecie w Zurichu. Przedstawione podejście do tematyki sieciowej integracji danych rozszerza zagadnienie optymalizacji ruchu o aspekt priorytetów wykonania. Propozycja jest nowatorska, w stosunku do obecnie stosowanych metod trasowania sieci.