

IB109 poznamky – 2020

01 - Uvod

Subeznost – existencia dvoch alebo viacerých procesov súčasne

Výpočetný paralelizmus

- Na všetkých úrovniach, od registrov až po koexistenciu rozsiahlych distribuovaných aplikácií
- Chceme ho kvôli vyššiemu výkonu
- Zároveň je nutný kvôli priestorovej distribúcii
- To však za cenu zložitejšieho návrhu a ceny

Sekvenčné algoritmy, složitostná trieda NC

Dovody pre vývoj paral. App.:

- Výkonnosť – efektívne využitie agregovanej výp. Prostriedky -> rýchlejší výpočet
- Provídateľnosť – v prípade kedy agregácia výp. Sily nie je voľba, ale nutnosť -> aby sa výpočet skončil v rozumnom čase, pretože kvôli náročnosti by inak trval príliš dlho
- Bezpečnosť – duplikácie časti systému, ak by niektorá havarovala alebo bola narušená jej dôvera
- Cena – agreg. Výp sila je lacnejšia a zároveň je lacnejšie udržiavať nesúvisiace časti aplikácie

Abstraktný model výpočetného systému:

Procesor – Datová cesta – Datové úložisko

Všetky časti systému môžu byť úzkym miestom voči výkonnosti aplikácie

Paralelizmus je prirodzený spôsob prekonania úzkych miest

Procesory

Neustála potreba zvyšovať výkon, pričom výkon je spájaný s **Moorovým zákonom**

Moorov zákon

- pomenovaný podľa spoluzakladateľa Intelu, Gordona Moore
- Hovorí že počet tranzistorov sa v procesoroch zdvojnásobi cca každých 18 mesiacov -> NEDAL sa úplne dodržať v jednom jadre, začalo sa pre zvyšovanie výkonu zavádzať viacero jadier, viz nissie

Metody zvyšovania výkonu procesoru

- Zvyšovanie frekvencie vnútorných hodín
- MULTIPLICITA, PARALELIZMUS

Pre pekne vyobrazenie moorovho zákona a počtu tranzistorov, vygoogli si "transistor counts and moore law"

Trendy vo vývoji procesorov

- Nedari sa zvyšovať výkon jedného jadra tak ako by sa chcelo
- Fyzikálne zákony bránia neustálej miniaturizácii – okolo 5 nm sa nedajú udržať elektróny v atóme
- Súčasne je to 14-16nm

Riesenie

- Viacero jadier procesoru
- Pravdepodobne takýto vývoj bude pokračovať aj do budúcnosti
- Je to odklon od monolitických jadier k desiatkam menších špecializovaných či hybridným riešeniam

Dosledok

- Sekvenčné algor. nemôžu naďalej rásť z rastúceho výkonu CPU
- Inu PARALELIZMUS VÝPOČTU JE NEVYHNUTNÝ SMER VÝVOJA

Datové cesty – paralelizmus v komunikácii

- Vetsia priepustnosť komunikačných liniek s následkom nižšej latencie
- Robustnosť a spoľahlivosť komunikačných liniek
- Napríklad sírka zbernice 32/64/128 bitov a dualband wifi router

Pamät

- Výkon procesorov prevyšuje výkon ostatných komponentov
- To znamená, že cesta procesor – pamäť – disk je príliš dlhá
- Doba pre získanie jednotky informácie z pamäte rastie s vzdialenosťou miesta uloženia od miesta spracovania

Viacurovnové uloženie informácií

Registry procesoru

L1/L2/L3 cache

Operačná pamäť

Cache I/O zariadení

Magnetické/optické disky

Cache

- Kopia casti dat v rychlo dostupnom mieste
- Moze a nemusí byt kontrolovana userom alebo OS

Priklad roznych moznosti kontroly

- L1/L2 cache v ramci CPU nemozu byt kontrolovane programatorom
- I/O efektívne algoritmy
 - o Obchadzaju virtualizáciu pameti kontrolovanou OS a miesto toho realizuju vlastny sposob pouzitia operacnej pamete ako cache pre data na disku

Multiplicita pametovych modulov

- Vacsie mnozstvo ulozitelnych dat
- Vacsie mnozstvo liniek do pameti, tj, vyssia priepustnost
- Vacsia rezie na udrzanie konzistencie
- Napr diskove pole, P2P siete, NUMA architektury
 - o NUMA – viac procesorove pocitace s viacerimi pametovymi modulmi usporiadanymi tak, ze pristup jedneho procesoru do roznych pametov modulov je rozne rychly. Obvykle je pamet blizko k urcitemu CPU, ale architektura umoznuje citanie danej pameti inym cpu

Paralelne vypocty

Paralelizmus z pohladu OS

Multitasking na jednom jadre

- Falosny multitasking, aplikacie sa na CPU striedaju pocas behu
- Zdanlive “bezi” viacero aplikacii
- Jednotka planovania OS je process

Multitasking na viacerych vypocetnych jadrach

- Rozne aplikacie maju priradaene rozne vyp. Jadra
- Inak standardny multitasking na kazdom jednom jadre
- Jednotka planovanie OS je process

Multitasking a multithreading

- Kazda aplikacia moze mat viac vypoc. Vlaken
- Vlakena sa v behu na jednom CPU jadre striedaju
- Vlakena jednej aplikacie mozu bezat na roznych jadrach
- Jednotka planovania OS je vlakno

Flynova klasifikacia

Single Instruction Single Data - SISD

- V daný okamžik je zpracována jen jedna instrukce a jedním datovým proudem
- Klasický sekvencný výpočet

Single Instruction Multiple Data - SIMD

- Jedna instrukce je vykonávána a více datovými proudy
- Velikost instrukce CPU, architektura GPU

Multiple Instruction Multiple Data - MIMD

- Nezávislý současný běh dvou a více SISD, SIMD přístupů

Multiple Instruction Single Data – MISD

- Prakticky se nevyskytuje

Distribuovaný/ paralelní systém

- Specifikovaný po částech, tj. procesech
- Chování systému vzniká interakcí současné bezcílové procesů
- Emergentně jeví, tj. chování, které není nutně explicitně naprogramováno/ plánováno, ale nastane shodou okolností

Synchronizace – omezení na překladání a průběh akcí jednotlivých procesů distribuovaného systému

Komunikace – přenos informací z jednoho procesu na druhý

Vybrané problémy distrib. Systémů

Javy

- Nekonzistentní vize konzistentního světa
- Vzájemná interferencí

Rizika

- Race-condition, tj. pokus o změnu dat dvěma procesy naráz, výsledné data závisí od pořadí vykonání procesů - nejasné
- nedeterministické chování
- Uváznutí (deadlock, livelock)
- Starnutí, hladovění (starving)
- Přetecení buferů, problém producent-konzument (producent nemůže produkovat, ak je bufer konzumenta plný, konzument nemůže konzumovat, ak je bufer prázdný)
- Zbytečná ztráta výkonu aktivním čekáním

Efektivní algoritmus se snaží sám zjistit, ak paměť je vhodná, například database.

8 CPU, 9 vlakien -> 1 vlakno sa prehadzuje medzi 8 cpu, ostatnych 8 vlakien pevne pridelenych, pretoze ak by boli zvlast 2 vlakna na jednom CPU, tak by bolo dane CPU 2x pomalsie

Dovody vyssiej narocnosti vyvoja paralelnych systemov

- Nutnost specifikace souběžných úkolů.
- Nutnost specifikace koordinace úkolů.
- Paralelní algoritmy.
- Nedostačující vývojová prostředí.
- Nedeterminismus při simulaci paralelních aplikací.
- Absence reálného modelu paralelního počítače.
- Rychlý vývoj a zastarávání použitých technologií.
- Výkon aplikace náchylný na změny v konfiguraci systémů

02 – programovanie v prostredi so zdielanou pametou

HW platformy

Paralelne systemy so zdielanou pametou

- Systemy s viac procesormi
- Systemy s viac jadernymi procesormi
- Systemy s procesormi so zabudovanim SMT – Simultanny multithreading
- Kombinacia

Rizika paralelnych vypoctov na sucasnych procesoroch

- Mnohe optimalizacie na urovni procesorov boli navrhnuté tak, aby zachovaly semantiku sekvencnych programov

POZOR na

- Preusporiadanie instrukcii
- Odloženie zapisu do pameti

SMT – Simultanny multithreading

Princip

- Procesor vyuziva prazdne cykly sposobene latenciou pameti k vykovananiu instrukcii ineho vlakna
- Vyzaduje duplikaciu urcitych casti procesoru, napr. Registrov
- Vlakna zdielaju cache
- Prepnutie vlakna urobi kopiu dat

Priklad: Intel Pentium 4

- Hyper-Threading Technology – **HTT**
- OS s podporou SMP(symetricky multiprocessing) vidi system s SMT/HTT ako viac jadrový system
- Az 30% narast vykonu, ale vzhľadom k sdielanej cache moze byt rychlost vypoctu jedneho vlakna nizsia

Viac-jaderne procesory – multicore

Viac plnohodnotných procesorov v jednom chipe

Vyhody

- Efektivnejšia cache koherencia na najnizsej urovni
- Nizsie naklady pre koncového usera

Nevyhody

- Viac jadier emituje vacsie zbytkove teplo(globalne oteplovanie)
- Takt jedneho jadra byva nizsi
- Automaticke docasne podtaktovanie/pretaktovanie
- Jadra zdielaju datovu cestu do pameti – pravdepodobne uzka cast, teda bottleneck

Realita

- Viac jadrove procesory so SMT
- Intel Core i7 – hexa core so SMT ma az 12 paralelných jednotiek

Paralelizmus v prostredí so sdielanou pametou

Idealizovaný model

- Tu sa riesi návrh paralelného algoritmu
- Jednotlive vypocetne jadra paralelného system pracuju uplne nezávisle
- Prístup k datám pameti je bezcasový a navzájom sa jednotlive pristupy vylucuju

- Komunikacia uloh prebieha atomicky cez zdielane datovy struktury

Realita

- Na tejto urovni programator riesi technicku implementáciu paral. Alg.
- Pristupy do pameti cez zbernicu su prilis pomale pre CPU
- Registry procesorov a cache pameti – rychle kope maleho mnozstva dat na roznych miestach datovej cesty
- Problem koherencie dat – data rovnake by mali byt rovnake pre vsetky procesy

Paralelne ulohy v kontexte OS – procesy a vlakna

Procesy

- Skryvaju pred ostatnymi procesmi svoje vypocetne prostriedky
- Pre riesenie paralelnej ulohy je potreba medzi procesova komunikacia – **IPC**
Zdielane pametovy segment, sokety, pomenovany a nepomenovane rury

Vlakna

- Existuju v kontexte jedneho procesu
- V ramci rodicovskeho procesu zdielaju vypocetne prostriedky
- Komunikacia prebieha cez zdielane datove struktury
- Ucelom interakcie je skor synchronizacia nez transport dat
- Subjekty procedury planovania

Priklad pouzitia vlakien vramci procesu

*

Vlakno

- Realizuje vypocet, tj sekvenciu instrukcii
- Kazd process je tvoreny aspon jednym vlaknom
- Hlavne vlakno procesu vytvara dalsie vlakna

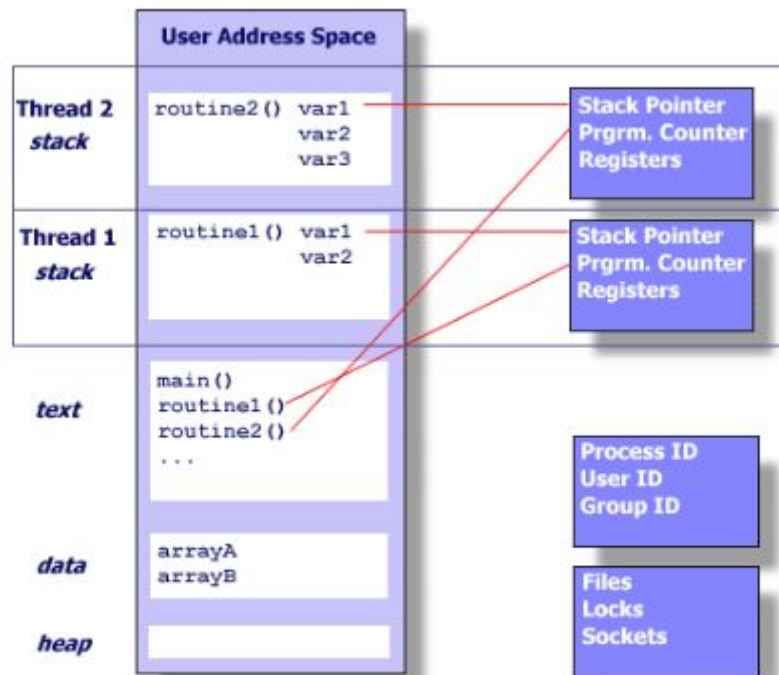
Proces ma

- Itentifikatory procesu a vlastnika
- Premenne prostredia, pracovny adresar
- Kod
- Registry, zasobnik, halda
- Odkazy na otvorene suboy a zdielane knihovne
- Reakcie na signaly
- Kanally IPC

Vlakna maju private

- Zasobnik

- Registry
- Frontu signalov



Argumenty pro pouzivanie vlakien miesto procesov

Vykon aplikace

- Vytvorenie procesu je vyrazne drahsie nez vytvorenie vlakna
- Zmena dat prevedena vramci jedneho vlakna je viditelna v kotnexte celeho procesu
- Komunikacia medzi vlaknami spociva v predavani referencie na data, nie predavaniu obsahu
- Predavanie referencie sa deje vramci jedneho procesu, operacny system nemusí riesit skryvanie dat a pristupove prava

Nevyhody

- Vlakna nemaju ziadne "sukromie"
- Zdielane globalne premenne

Efektivne vyuzite cache

Princip cache

- Mensia, ale rychlejsia pamet na pomalej datovej ceste
- Pri prvom nacistani sa okolie citanej informacie ulozi do cache
- Pri naslednom citani z okolia povodnej informacie staci citat z cache

Koherence

- Soulad dat ulozenych v pameti pocitaca, respective cache
- Je zajistene, ze existuje prave jedna platna hodnota asociovana s danym pametovym miestom
- Z dovodu rychlosti su zapisy procesoru do pameti odkladane a zdruzovane, blizsie specifikacie chovanie procesoru v tomto ohlade je dana **pametovym modelom** daneho CPU

Suvisiace pojmy

- Cache line – atomicky pametovy blok ulozeny v cache
- “hit ratio” - cislo vyjadrujuce uspesnost obsluzenia poziadavkov na data datami ulozenymi v cache
- “vyliatie cache” – procedura aktualizacie dat v pameti hodnotami ulozenymi v cache

Zasady efektivneho pouzitia cache

- Casova lokalita – pristupy v malom casovom intervale
- Priestorova lokalita – pristup k datom ulozenym adresne blizko seba
- Zarovnana alokacia pameti – napr. Memalign v GNU C

False sharing

- Paralelny cache koherentny system s viac procesormi
- Program s niekolkonasobne viacero vlaknami
- Pole hodnot `int pole[nr_of_threads]`
- Vlakna pocitaju vyslednu hodnotu typu `int` a k vypoctu si ukladaju medzivysledok typu `int`

Varianty implementacie

- A) Kazde vlakno opakovane zapisuje do datoveho pola integerov na poziciu urcenu jeho ID
- B) Kazde vlakno zapisuje do lokalnej premennej a pred skoncenim nakopiruje hodnotu do pola na poziciu urcenu jeho ID

*

Pouziva sa cache line filler na zvyshenie vykonu.

Pricom nie je garantovane ze vsetky informacie ulozene do explicitne definovanych premennych budu zapisane do pameti a zaroven premenne mozu byt realizovane registrom procesoru.

A ako dosledok, postupnost hodnot ulozenych do jednej premennej vramci jedneho vlakna moze byt videna **ciastocne** alebo **vobec**.

Udrzovanie koherencie cache pameti je naklade a zaroven sucasne prekladace **nevynucuju** aby kazda vypocitana hodnota bola ulozena do pameti.

Ako dosledok, modifikacie v jednom vlakne CPU, sa nemusi prejavit v inom vlakne na inom CPU.

Problem:

- Vlakno P0 neskonci alebo nezaznamena zmenu premennej x
- Moze zavisiet aj na stupni optimizacie prekladacu
 - o Chyba sa neprejaví pri g++ -O0
 - o Ale prejaví sa pri g++ -O2

Riesenie

- Je nutne oznacit premenu ako nestalu premenu, tj, **volatile**(c, c++)
- Prekladac zaisti aby dana premenna nebula realizovana iba na urovni registrov CPU, ale pred a po kazdom pouziti bola nacistana/ulozena do pamete

Rozlisujeme rozdielne typy volatile:

- Umiestnenie pred alebo za datovy typ v definicii premenej
- Rozlisujeme:
 - o Ukazatel na nestalu premennu
 - `Volatile T *a`
 - o Nestaly ukazatel na nestalu premennu
 - `Volatile T * volatile a`
 - o Nestalu premennu
 - `Volatile T a`
 - `T volatile a`

Pripady kedy je nutne pouzit volatile:

- Premenna zdielana medzi subeznymi vlaknami/procesmi
- Premenna zastupujuca vstupny/vystupny port
- Premenna modifikovana procedurou obsluhujucu prerusenie

Presny popis toho ako processor manipuluje so zapismi dat do pameti je sucastou specifikacie procesoru, jedna sa o tzv pametovi model.

*Napriklad, zapisy hodnot do pamete niesu vykonane v okamihu zpracovania instrukcie, ale su odkladane na neskor a zhukovane. Preto napriklad pocitanie do 200000 na systeme s jednou vykonovou jednotkou je vzdy 200000, ale na paralelnych architekturah je castokrat mensie.

Pozor:

- Poradie zapisov do pameti podľa instrukcii CPU sa nemusí zhodovať so skutočným poradím zapisu hodnôt do pameti
- Pametový model garantuje korektnosť **iba** pre sekvencne programy

*Paralelné programovanie si vyžaduje opony bod na úrovni HW, inak takmer nemožné.

Pametová bariera

- HW primitívum pre synchronizáciu stavu pamäte a stavu procesorov v danom mieste programu
- Na súčasných procesoroch realizované instrukciou **mfence**
- Na hardwarovej úrovni prevedie serializáciu všetkých *load* a *store* instrukcií, ktoré sa vyskytujú pred instrukciou *mfence*. Táto serializácia zaisťuje, že efekt všetkých instrukcií pred instrukciou *mfence* bude globálne viditeľný pre všetky instrukcie nasledujúcich za instrukciou *mfence*

HW podpora pre atomicke instrukcie

- Pametová bariera nerieši problém atomických instrukcií ako sú napríklad TEST-AND-SET, COMPARE-AND-SWAP
- Instrukcie spomenutého typu sú však pre účely efektívneho paral. progr. veľmi vhodné

Dalsia HW podpora

- Alpha, Mips, PowerPC, ARM: instrukcie typu LL(Load Linked)/SC(Store Conditional)
 - o LL vráti súčasnú hodnotu pamäte na adrese, pričom nasledujúci SC uloží na rovnakú adresu nové data, iba ak nedošlo k zmene na adrese od času vykonania LL
- X86 architektúra
 - o Lock – nasledujúca instrukcia, ktorá zapisuje do pamäte prebehne atomicky, a jej efekt bude hneď globálne viditeľný
 - o XCHG – prehodí obsah registru a pametového miesta -> obsahuje z definície prefix *lock*

Možnosti realizácie atomických instrukcií na úrovni kodu

1. Jazyk symbolických adries – JSA -> fancy pojem pre assembler
2. Zabudované funkcie prekladáča, od GCC 4.1
 - a. `type __sync_val_compare_and_swap`
 - b. `type __sync_fetch_and_add`
3. súčast programovacieho jazyka – napríklad C++ rev 11, Java, ...

*

03 – programovanie v prostredí so zdieľanou pametou vol 2.

Rizika spojene so zdielanou pametou

Paralelne programy mozu byt opakovanom spusteni zdanlive nahodne vykazovat rozne chovania

Vysledok programu moze zavistiet na absolutnom poradí prevedenia instrukcii programu, tj preloženie instrukcii zúčastnených procesov/vlákien

Race condition

- nedokonalosť paralelného programu, ktorá sa prejavuje hore spomínaným spôsobom sa označuje **race condition**, zkratkene **race**

Atomicita operácii

- jednoduchý príkaz vo vyššom programovacom jazyku neodpovedá nutne jednej instrukcii procesoru
- v moderných OS je každé vlákno podrobené plánovaciemu procesu
- vykonanie postupnosti instrukcii procesoru odpovedajúceho jednému príkazu vyššieho programu. Jazyk môže byť prerušenie a preloženie vykonaním instrukcie iného vlákna
- napríklad, priradenie čísla do premennej efektívne môže znamenať nabitie do registru, prevedenie aritmetickej operácie a uloženie výsledku do pamäte
 - o * môže sa stať, že sa efekt jedného priradenia do zdieľanej globalnej premennej úplne stratí

Relatívna rýchlosť výpočtu *

Neda sa spoľiehať na súčasný subeh vlákien a teda aj na relatívnu rýchlosť výpočtu jednotlivých vlákien

Uviaznutie – Deadlock

- Pokiaľ majú vlákna **neusporiadané inkrementálne** požiadavky na **unikátne** zdieľané zdroje, môže dôjsť k deadlocku

Ak uviaznu 2 vlákna z viacerých, jedna sa stále o uviaznutie.

Čiastočný deadlock sa veľmi náročne detekuje

Hladovanie, Starnutie, nepogrese – Livelock *

Jedná sa o jav, kedy aspoň jedno vlákno nie je schopné kvôli pohľadom k paralelnému subehu s iným vláknom pokročiť vo výpočte za danú hranicu

Thread-safe procedura

Označuje procedúru či program, ktorého kód je bezpečne prístupný (vzhľadom k semantike výstupu a stabilite výpočtu) subezne s niekoľkými vláknami bez nutnosti vzájomnej synchronizácie

Ak sa k datovym strukturam pristupuje naraz z viacerych vlakien, nejedna sa o thread-safe pripad

Knihovnové funkcie nemusia byt thread-safe

- Rand() -> rand_r()

Re-entrantna procedura

Jedna sa o procedure, ktorej prevedenie moze byt vramci jedneho vlakna prerusene, kod kompletne vykonany od zaciatku do konca v ramci tej istej ulohy, a potom znovu obnoveny/dokonceny po preruseni kodu

Je to termin z cias ked este neboli multitaskingove OS

Re-entrantna procedura nemusí byt thread-safe a rovnako thread-safe procedura nemusí byt re-entrantna *

Nebezpecne akcie vzhľadom k paralelnemu zpracovaniu

- Nekontrolovany pristup k global. Premen. A halde
- Uchovavani stavu procedury do globalnych premennych
- Alokacie a dealokacie zdrojov globalneho rozsahu(subory,...)
- Nepriamy pristup k datam skrz odkazy alebo ukazatele
- Viditelny vedlajsi efekt - > napr modifikacie nestalych premennych

Bezpecne strategie

- Pristup iba k lokalnym premennym(zasobnik)
- Kod je zavisly iba na argumentoch danej funkcie
- Akedkoľvek volane podprocedury a funkcie su thread-safe *

Pristup k sdíelaným premenným

Pristupovanie k sdíelaným premenným je nekorektne

Všetky modifikácie a neatomicke citania globalnych premennych musia byt **serializovane**

Kritická sekcia

- Cast kodu ,ktorej prevedenie je neprelozitelne instrukciami ineho vlakna
- Realizacia kritickej sekcie musi byt odolna voci planovaniu

Zamykanie

- Vlakno vstupujuce do volnej kritickej sekcie svojim vstupom znemozni pristup ostatnym vlaknam
- Ostatne vlakna cakaju pred vstupom do kritickej sekcie
- Pri odchode z kritickej sekcie sa zamok uvolni a ziska ho **nahodne** cakajuce vlakno

Jednoduché riešenie zamku

- Zdieľaná atomická prístupová bitová premenná, ktorej hodnota indikuje prítomnosť procesu či vlákna v priradenej kritической sekcii
- Premenná je manipulovaná pri vstupe a výstupe z/do kritической sekcii
- Vyžaduje podporu HW pre atomickú manipuláciu

Aktivné čakanie – spinlock

Dokiaľ neúspeje, vlákno sa bude snažiť dostať do sekcii opakované

Uspávanie

Procesy/vlákná sa po neúspechu vstúpiť do kritической sekcii uspia

Zobudia sa po vypršaní časového limitu alebo explicitne iným vláknom

Riziká zamykania

- Uviaznutie, starnutie, zníženie výkonu

Manipulácia so zamkom vyvoláva vytláčanie cache pamäte

Mnoho prístupov k zamykaným premenným môže byť úzkym miestom výkonu aplikácie, z princípu sa ale nedať odstrániť

Petersonov algoritmus – spinlock, user-space

- Spravodlivý algoritmus priradenia vzájomného vylúčenia
- Nesposobuje starnutie ani uviaznutie
- Vyžaduje atomicke zápisy do premenných
- Citlivý na vykonávanie inštrukcií mimo poradia

POSIX thread API

Základné delenie funkcionality

Správa vlákien

- Vytváranie, oddelovanie a spojovanie vlákien
- Funkcie na nastavenie a zistenie stavu vlákna

Vzájomné vylúčenie – mutexes

- Vytvorenie, nariadenie, zamykanie a odomýkanie mutexov
- Funkcie na nastavenie a zistenie atribútov spojených s mutexami

Podmienkové/podmienenné premenné – conditional variable

- Služi pre komunikáciu/synchronizáciu vlákien
- Funkcie na vytváranie, nariadenie, “čakanie na” a “signalizovanie pri” specifickej hodnote podmienkovej premennej
- Funkcie na nastavenie a zistenie atribútov premenných

Posix standard

Cez 60 API funkcii

- #include <pthread.h>
- Preklad s volbou -pthread

Mnemotechnicke predpony funkcii

- pthread_, pthread_attr_
- pthread_mutex_, pthread_mutexattr_
- pthread_cond_, pthread_condattr_
- pthread_key_

Pracuje sa so skrytmi objektmi – opaque objects

- Objekty v pameti o ich podobe programator nic nevie
- Pristupovane k vyhradne pomocou odkazu (handle)
- Nedostupne objekty a neplatne(dangling) referencie

Atributy objektu

Idea

- Vlastnosti vsetkych vlakien, mutexov a podmienkovych premennych nastavovane specialnymi objektmi
- Niektore vlastnosti entity musia byt specifikovane uz v dobe vzniku entity

Typy atributovych objektov

- Vlakna: pthread_attr_t
- Mutexy: pthread_mutexattr_t
- Podmienkove premenne pthread_condattr_t

Vznik a destrukcia

- Funkcie _init a _destroy s odpovedajucou predponou
- Parameter odkaz na odpovedajuci atributovy object

Sprava vlakna

Vytvorenie vlakna

- Kazdy program ma jedno hlavne vlakno
- Dalsie vlakna musia byt explicitne vytvorene programom
- Kazde vlakno (i vytvorene) moze dalej vytvarat dalsie vlakna
- Vlakno je vytvorene funkciou pthread_create
- Vytvorene vlakno je hned pripraveno k pouzitiu
- Moze byt planovacou spustene skor nez sa dokonci volanie vytvaracej funkcie
- Vsetky data potrebne pri spusteni vlakna musia byt pripraveno pred volanim vytvaracej funkcie
- Maximalny pocet vlakien je zavisly na implementacii

*

pthread_create argumenty

thread_handle je odkaz na vytvorene vlakno

attribute odkaz na atributy vytvoreneho vlakna(NULL pre prednastavene nastavenie atributov)

thread_function ukazatel na funkciu noveho vlakna

arg ukazatel na parametry funkcie thread_function

pri uspesnom vytvorení vlakna vracia 0

Ukoncenie vlakna nastava

- Volanim funkcie pthread_exit
- Pokial skonci hlavna funkcia rodicovskeho vlakna inak nez volanim pthread_exit
- Ak je zrusene inym vlaknom pomocou pthread_cancel
- Rodicovsky process je ukonceny (nasilne alebo volanim exit)

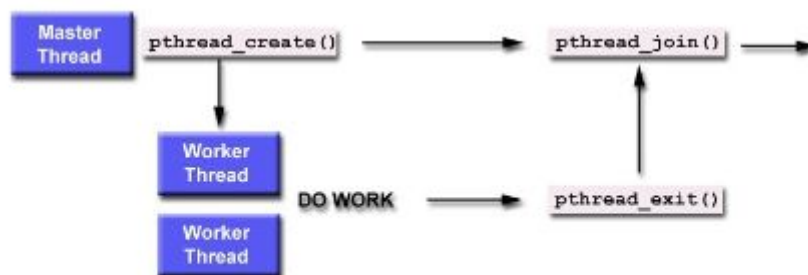
Void pthread_exit (void * value)

- Ukoncuje beh vlakna
- Odkazy na prostredie procesu (subory, IPC, mutexy) otvorene vramci vlakna sa nezatvaraju
- Data patriace vlaknu musia byt uvolnene pred ukoncenim vlakna(system prevedie uvolnenie prostriedkov az po skonceni rodicovskeho procesu)
- Ukazatel *value* predany pri spojeni vlakien

*

Int pthread_join (pthread_t thread_handle, void **ptr_value);

- Caka na dokoncenie vlakna *thread_handle*
- Hodnota *ptr_value* je ukazatel na pointer specifikovany vlaknom *thread_handle* pri volani *pthread_exit*



-
- Nutny napríklad pokiaľ main ma vracať zmysluplnú návratovú hodnotu

Vlastnosti vlakien

Nespojitelne vlakna – detached threads

- Nemozu byt spojené volanim funkcie pthread_join
- Setria systemove prostriedky
- Prednastavene nastavenie typu vlakna nie je vzdy zrejmé, preto je doporučene typ vlakna explicitne nastavit

*

Velkost zasobniku

- Minimalna velkost zasobniku nie je urcena
- Pri velkom pocte vlakien sa casto stave ze vyhradene miesto pre zasobnik je vycerpane
- Posix umoznuje zistit a nastavit poziciu a velkost miesta vyhradeneho pre zasobnik jedneho vlakna

*

Zrusenie vlakna

Int pthread_cancel (pthread_t *thread_handle);

- Vysle ziadost o zusenie vlakna *thread_handle*
- Adresovane vlakno sa moze a nemusí ukončit
- Vlakno moze ukončit sameho seba
- Pri zrusení sa prevedie vycistenie dát spojených s rúseným vlaknom
- Funkcia skončí po odeslaní ziadosti (je noblokujuca)
- Navratovy kod 0 znaci, ze adresove vlakno existuje, nie ze bolo/bude zrusene

Funkcia *pthread_cleanup_push* registruje cistiaci kod

Pomocou *pthread_setcancelstate* sa dá nastavit ci je vlakno zrusitelne a ako

Vzajomne vylucenie

Motivacia

- Viacero vlakien vykona nasledujuci kod
If (my_cost < best_cost) best_cost = my_cost;
- Nedeterministicky vysledok pre 2 vlakna a hodnoty:
Best_cost==100, mycost@1==50, my_cost@2==75

Riesenie

- Umiestnenie kodu do kritickej sekcie
- Pthread_mutex_t

Inicializacia mutexu

```
int pthread_mutex_init (pthread_mutex_t  
*mutex_lock, pthread_mutexattr_t *attribute)
```

- Parameter *attribute* specifikuje vlastnost zamku
- NULL znamena defaultne nastavenie

```
int pthread_mutex_lock (pthread_mutex_t *mutex_lock)
```

- Volanie tejto funkcie zamyka *mutex_lock*
- Volanie je blokujuce, dokial sa nepodari mutex zamknut
- Zamknut mutex sa podari iba raz jednemu vlaknu
- Vlakno, ktoremu sa podari mutex zamknut je v kritickej sekcii
- Pri opusteni kritickej sekcie, je vlakno "povinne" mutex odomknut
- Az po odomknuti je mozne mutex znova zamknut
- Kod vykonany vramci kritickej sekcie je po odomknuto mutexu globalne viditelny(pametova bariera)

```
int pthread_mutex_unlock (pthread_mutex_t *mutex_lock)
```

- Odomyka *mutex_lock*

Pozorovanie

- Velke kriticke sekcie znizuju vykon aplikacie
- Prilis casus a travi v blokujucom volani funkcie *pthread_mutex_lock*

Riesenie:

```
int pthread_mutex_trylock (pthread_mutex_t *mutex_lock)
```

- Pokusi sa zamknut mutex
- Vpripade uspechu vrati 0
- Vpripade neuspechu vrati *EBUSY*
- Zmysluplne vyuzitie komplikuje navrh programu
- Impementacia byva rychlejsia nez *pthread_mutex_lock*(nemusi sa manipulovat s frontami cakajucich procesov)
- Ma zmysel aktivne cakat opakovanym volanim trylock

Vlastnosti(atributy) mutexov

Normalny mutex

- Iba jedno vlakno moze jedenkrat zamknut mutex
- Pokial sa vlakno, ktora ma zamknuty mutex pokusi znova zamknut rovnaky mutex, dojde k uvaznutiu

Rekurzivny mutex

- Dovoľuje jednemu vlaknu zamknut mutex opakovane

- K mutexu je asociovaný citac (pocítadlo) zamknutí
- Nenulový citac znamená zamknutý mutex
- Po odemknutí je nutné zavolat unlock tolikrát, kolikrát bylo voláno lock

Normalný mutex s kontrolou chyby

- Chová se jako normální mutex, akurát při pokuse o druhé zamknutí ohlásí chybu
- Pomalejší, typicky používané dočasně počas vývoje aplikací, potom nahrazeny normálním mutexem

```
int pthread_mutexattr_settype_np (pthread_mutexattr_t
*attribute, int type)
```

- Nastavení typu mutexu
- Typ určený hodnotou proměnné *type*
- Hodnota *type* může být:
 - PTHREAD_MUTEX_NORMAL_NP
 - PTHREAD_MUTEX_RECURSIVE_NP
 - PTHREAD_MUTEX_ERRORCHECK_NP

04 – POSIX THREAD CONT., WIN32 THREADS

Základné delenie:

Správa vláken

- Vytváření, oddělování a spojování vláken
- Funkce na nastavení a zistenie stavu vlákna

Vzájomne vylucenia – mutexes

- Vytváření, nicení, zamykání a odomykání mutexů
- Funkce na nastavení a zistenie atributů spojených s mutexami

Podmienkove/podmienene premenne – conditional variable

- Slouží pro komunikaci/synchronizaci vláken
- Funkce na vytváření, nicení, “čekaní na” a “signalizování při” specifické hodnotě podmínkové proměnné
- Funkce na nastavení a zistenie atributů proměnných

Podmienkove premenne

Motivace I

- Často na jednu kritickou sekci čeká více vláken
- Aktivní čekání – permanentní zátěž CPU

- Uspavanie timeoutom – netrivialna rezia, omedzena frekvencia dotazovania sa na moznost vstupu do kritickej sekcie

Motivacia II

- Vlakno realizuje ucelenu, logicky oddelenu funkcionalitu
- Ta neni potrebna za celu dobu behu aplikacie
- Programatorom riadena docasna deaktivacia vlakna

Obene riesenia:

- Uspanie vlakna, pokiaľ vlakno ma/musi cakat
- Vzbudenie vlakna v okamihu ked je mozne pokracovat

Realizacia v POSIX Threads

- Mechanizmus oznacovany ako *podmienkove premenne*
- Podmienkova premena vyzaduje pouziti mutexu
- Po ziskani mutexu sa vlakno moze docasne vzdac tohoto mutexu a uspat sa (vramci danej podmienkovej premennej)
- Prebudenie musi byt explicitne signalizovane inym vlaknom

```
int pthread_cond_init (pthread_cond_t
*cond, pthread_cond_attr_t *attr)
```

- Inicializuje podmienkovu premennu
- Ak ma *attr* hodnotu NULL, pouzije sa defaultne chovanie

```
int pthread_cond_destroy (pthread_cond_t *cond)
```

- Nici nepouzivanu podm. Premennu a suvisiace datove struktury

```
int pthread_cond_wait (pthread_cond_t
*cond, pthread_mutex_t *mutex_lock)
```

- Uvolni mutex *mutex_lock* a zablokuje vlakno v spojeni s podmienkovou premennou *cond*
- Po navrate vlakno opet vlastni mutex *mutex_lock*
- Pred pouzitim musi byt *mutex_lock* inicializovany a zamknuty volajucim vlaknom

```
int pthread_cond_signal (pthread_cond_t *cond)
```

- Volane z ineho vlakna
- Signalizuje prebudenie jedneho z vlakien, uspanych and podmienkovou premennou *cond*

```
int pthread_cond_broadcast (pthread_cond_t *cond)
```

- Signalizuje prebudenie vsetkym vlaknam cakajucim nad podm. Premennou *cond*
-

```
int pthread_cond_timedwait (pthread_cond_t  
*cond, pthread_mutex_t *mutex_lock, const struct timespec  
*abstime)
```

- Vlakno bud zobudene signalom, alebo zobudene po uplynuti castu specifikovanom v *abstime*
- Pri vzbudení z dovodu uplynutia casu, vracia chybu *ETIMEDOUT*, a neimplikuje znovu získanie *mutex_lock*

Pouzitia podmienkových premenných *

Dalsie funkcie v POSIX threads

Globalne premenne specificke pre vlakno

Problem

- Vzhľadom k požiadavkom vytvárania reentrantných a thread-safe funkcií sa programátorom zakazuje používať globálne data
- Pripadne použitie globálnych premenných musí byť bezstavové a vykonané v kritickej sekcii
- Kladie obmedzenia na programátorov

Riesenie

- Thread specific data – **TSD**
- Globalne premenne, ktoré môžu mať pre každé vlákno inú hodnotu

Implementacia TSD

Standarne riesenie

- Pole idnexované jednoznačným identifikátorom vlákna
- Vlákna musia mať rozumne veľké identifikatory
- Snadný prístup k dátam patriacim iným vlaknám – potenciálne riziko nekorektneho kodu
- Výkonnostný problém a sanca false sharing ak sa do pola s premennymi pristupuje príliš často

Riesenie POSIX standard

- Identifikátor(kľuč) a asociovaná hodnota
- Identifikátor je globálny, asociovaná hodnota je lokálna premenná
- Kľuč – *pthread_key_t*

- Asociovaná hodnota – univerzálny ukazateľ, tj void *

Posix klucy

```
int pthread_key_create (pthread_key_t *key, void
(*destructor) (void*))
```

- Vytvorí nový kľúč – jedna sa o globálnu premennú
- Hodnota asociovaného ukazateľa je nastavená na NULL pre všetky vlákna
- Parameter *destructor* – funkcia, ktorá bude nad asociovanou hodnotou vlákna volaná v okamihu ukončenia vlákna, pokiaľ bude asociovaná hodnota nenulový ukazateľ
- Parameter destructor je nepovinný, lze nahradiť NULL

Použitie posix kľúčov

Zníženie kľúča a asociovaných ukazateľov

- Int pthread_key_delete(pthread_key_t key)
- Nevola žiadne *destructor* funkcie
- Programátor je zodpovedný za dealokáciu objektov pred zničením kľúča
- Funkcia destructor sa volá len keď končí vlákno a pthread_key_delete nie je null

Funkcie na získanie a nastavenie hodnoty asociovaného ukazateľa

- Void * pthread_getspecific (pthread_key_t key)
- int pthread_setspecific (pthread_key_t key, const void *value)

Rozne

Pthread_t pthread_self ()

- vráti unikátny systémový identifikátor vlákna – sama seba

int pthread_equal(pthread_t thread1, pthread_t thread2)

- vráti nenula pri identite vlákien thread1 a thread2

```
pthread_once_t once_control = PTHREAD_ONCE_INIT;int
pthread_once(pthread_once_t *once_control, void
(*init_routine) (void));
```

- prvé volanie tejto funkcie z akéhokolvek vlákna spôsobí vykonanie kódu *init_routine*. Ďalšie volania nemajú žiadny efekt
- init_routine môže byť napríklad inicializovanie premenných na začiatku atd

Ďalšie funkcie v POSIX threads

Planovanie(scheduling) vlákien

- není definované, väčšinou je defaultná politika dostávajúca

- POSIX threads poskytuje funkcie na definíciu vlastnej politiky a priorit vlákien
- Není požadovaná implementácia tejto časti API
- Chovanie nie je garantované

Správa priorit mutexu

Zdieľanie podmienkových premenných medzi procesmi

Vlákna a obsluha posix signalov

Read-Write zamky

Typické konštrukcie

Ctenari a pisari – WRRM mapa

Specifikácia problému

- Vlákna často citajú zdieľanú hodnotu, ktorá je relatívne menej často modifikovaná (Write-Rarely-Read-Many)
- Je chcené aby čítanie hodnoty mohlo prebiehať subezne

Možné problémy

- Subezný prístup dvoch vlákien-pisateľov môže vyústiť v nekonzistentné dáta alebo mať nežiaduci vedľajší efekt, napríklad memory leak
- Subezný prístup vlákna-pisateľa v okamihu čítania hodnoty iným vláknom-čtenárom môže vyústiť v čítanie neplatných, nekonzistentných dát

Ctenari a pisari – riešenie

Riešenie s použitím POSIX Threads

- Čítanie a modifikácia dát bude prebiehať v kritickej sekcii
- Prístup do kritickej sekcie bude riadený pomocou funkcie pthread_*

Dalšie požiadavky

- Vlákno-čtenár môže vstúpiť do kritickej sekcie pokiaľ v nej nie je alebo na ňu nečeka vlákno-pisateľ
- Vlákno-čtenár môže vstúpiť do kritickej sekcie pokiaľ sú v nej iné vlákna-čtenári
- Prístup vlákien-pisateľa je serializovaný a majú prednosť pred prístupom vlákien-čtenárov

Ctenari a pisari – implementácia

Jednoduché riešenie

- Použiť jeden pthread_mutex_t pre kritickú sekciu
- Vylúčiť subezný prístup vlákien-čtenára

Lepšie riešenie

- Implementuje nový typ zámku – `rwlock_t`
- Funkcia pracujúca s novým zámkom
 - `Rwlock_rlock(rwlock_t *l)` – vstup vlákna-čtenare
 - `Rwlock_wlock(rwlock_t *l)` – vstup vlákna-pisare
 - `Rwlock_ulock(rwlock_t *l)` – opustení libovolným vláknem
- Funkcie `rwlock` implementované s použitím podmienkových premenných z POSIX thread API

*

Barriery

Specifikace problem

- Synchronizační primitivum
- Vlákno je dovolené pokračovať po bariere až keď ostatné vlákna dosiahli bariery
- Naivná implementácia cez mutexy vyžaduje aktívne čakanie – nemusí byť vždy efektívne

Lepsie riesenie

- Implementácia bariery použitím podmienkovej premenej a počítadla
- Každé vlákno ktoré dosiahne bariery zvýši počítadlo
- Pokiaľ nie je dosiahnutý počet vlákien, podmienené čakanie

*

Problem

- Po dosiahnutí bariery všetkými vláknami, je mutex `count_lock` postupne zamknutý pre všetky vlákna
- Dolný odhad na dobu behu bariery je teda $O(n)$, kde n je počet vlákien participujúcich na bariere

Mozne riesenie

- Implementácia bariery metódou binárneho poľenia
- Teoretický dolný odhad na barieru je $O(n/p + \log p)$, kde p je počet procesorov

Chyby, okrem nezamykaného prístupu ku globalnej premennej

Typicke chyby – situace 1

- Vlákno V1 vytvára vlákno V2
- V2 požaduje data od V1
- V1 plní data až po vytvorení V2
- V2 použije neinicializované data

Typicke chyby – situace 2

- Vlákno V1 vytvára Vlákno V2
- V1 preda V2 pointer na lokálne data V1

- V2 pristupuje k datam asynchrone
- V2 pouzije data, ktora uz neexistuju(V1 skoncil)

Typicke chyby – situace 3

- V1 ma vyssiu prioritu nez V2, citaju rovnake ata
- Neni garantovane ze V1 pristupuje k datam pred V2
- Pokial ma V2 destruktivne citanie, V1 pouzije neplatne data

Ladenie programov s POSIX vlakny

Valgrind

- Simulacia behu programu
- Analyza jedneho behu programu

Nastroje valgrindu

- Memcheck – detektce nekonzistentneho pouzitia pameti
- Callgrind – jednoduchy profiler
 - o Kcachegrind -vizualizace
- Helgrind – detekce nezamykanych pristupov ku zdielanym premennym v POSIX thread programoch

Rozsirenie POSIX threads – nepovinna dle standard -> bariery, RW zamky

Dalsie sposoby synchronizacie

Synchronizacie procesov

Problem – jak synchronizovat procesy

- Mutexy z POSIX threads dle standard sluzia iba pre synchronizaciu vramci procesu
- Pre realizaciu kritickych sekcii v roznych procesoch je treba aj inych synchronizacnych primitive
- Podpora zo strany OS

Semaforey

- Citace pouzivane ku kontrole pristupov ku zdielanym zdrojom
- POSIX semaforey(vramci procesu)
- System V semaforey(medzi procesy)
- Daju sa vyuzit i pre synchronizaciu vlakien

Semafor

- Celociselny, nezaporny citac, ktoreho hodnota indikuje “obsadenost” zdielaneho zdroja
 - o Nula – zdroj vyuzivan a neni k dispozici
 - o Nenula – zdroj nevyuzivan, je k dispozici

- Sem_init() – inicializuje citac zadanou defaultnou hodnotou
- Sem_Wait() – znizi citac, pokiaľ moze a skonci, inak blokuje
- Sem_post() – zvyši citac o 1, prípadne zobudi čakajúce vlakno

Semafore vs mutexy

- Mutex moze odomknut iba to vlakno, ktore ho zamklo
- Semafor moze byt spravovany/manipulovany roznyimi vlaknami

Monitor

- Synchronizacne primitivum vyssie programovacieho jazyka
- Označenie kodu, ktorý moze byt subezne vykovany najviac jednym vlaknom
- JAVA – klucove slovo *synchronized*

Semafore, mutexy a monitory

- So semaformi a mutexmi je explicitne manipulovane programatorom
- Vzajomne vylucenie realizovane monitorom je implicitne, tj, explicitne zamykanie skrz explicitne primitive doplni prekladac

Vlakna v MS Windows

Vyssi programovaci jazyk

- C++11
- JAVA
- ...

Posix Thread pre Windows

- Existuje knihovna poskytujúca POSIX thread interface

Nativne rozhranie MS Windows

- Priame systemove volania – súčasť jadra OS
- Iba ramcovo podobna funkcionalita ako POSIX threads – tie iste veci riesene castokrat inak
- Na rozdiel od POSIX threads nema nepovinne casti (teda neexistuju rozne implementacie toho isteho)

Windows vs. POSIX Threads – Funkce

Windows	POSIX Threads
Pouze jeden typ HANDLE	Každý objekt má svůj vlastní typ (např. pthread_t, pthread_mutex_t, ...)
Jedna funkce pro jednu činnost. (např. WaitForSingleObject)	Různé funkce pro manipulaci s různými objekty a jejich atributy.
Typově jednodušší rozhraní (bez typové kontroly), čitelnost závislá na jménech proměnných.	Typová kontrola parametrů funkcí, lepší čitelnost kódu.

Win32 vs. POSIX Threads – Synchronizace

Windows	POSIX Threads
	Mutexy Podmínkové proměnné Semaforey
Signalizace pomocí událostí.	Signalizace v rámci podmínkových proměnných.
Jakmile je událost signalizována, zůstává v tomto stavu tak dlouho, dokud ji někdo voláním odpovídající funkce nepřepne do nesignalizovaného stavu.	Signál zachycen čekajícím vláknem, pokud takové není, je signál zahozen.

Win32 vs. POSIX Threads – Základní mapování

Windows	POSIX Threads
CreateThread	pthread_create pthread_attr_*
ThreadExit	pthread_exit
WaitForSingleObject	pthread_join pthread_attr_setdetachstate pthread_detach
SetPriorityClass SetThreadClass	Pouze nepřímo mapovatelné. setpriority sched_setscheduler sched_setparam pthread_setschedparam ...

Win32 vs. Linux/UNIX – Mapování synchronizace

Windows	Linux threads	Linux processes
Mutex	PThread Mutex	System V semafor
Kritická sekce	PThread Mutex	—
Semafor	PThread podm. proměnné POSIX semafor	System V semafor
Událost	PThread podm. proměnné POSIX semafor	System V semafor

Win32 vs UNIX – pozorovanie

Pozice vláken v MS Windows

- Silnejšie postavenie než vlákna v Linuxe
- Synchronizačné prostriedky fungujú i medzi procesmi
- Vlákna vo vláknach (Processes-Threads-Fibers)
 - o Scheduling Fibers nerieši OS, ale kooperatívne – vlákna sa sami rozhodujú či sú na rade a podobne, chyba externých zásahov OS

- User-mode-scheduling(UMS) – kooperativne planovanie

Vyhodny len jedneho typu zastupujuceho objekty

- Jedinou funkciou sa da cakat na nekonkretne vlakno
- Jednou funkciou sa da cakat na vlakno a na mutex

05 – implementacia lock-free datovych struktur

Snazime sa vyhnut pouzitiu zamykacich primitive pre zvysenie vykonu

Klasicka skola viacvlaknoveho programovania

- Pristup k zdielanym datam musi byt chrany
- Pristupy k datam sa musia serializovat s vyuzitim roznch synchronizacnych primitive(mutexy,semafony,monitory_
- Vlakna operuju s datami tak, aby sa tieto operacie javili ostatnym vlaknam ako atomicke operacie

Problemy

- Omeskanie pri pristupu k zdielanym datam
- Uviaznutie, zivost, ferovost
- Korektnost implementacie
 - o Tazko dokazat a tazko nevieme povedat na vysej urovni ci sa jedna o atomicku operaciu
- Atomicita operacii

Lock-free programovanie

- Programovanie paralelnych aplikacii bez pouzitia zamykania alebo inych makro-synchronizacnych mechanizmov

Vlastnosti lock-free programovania

- **Vyzaduje HW podporu**
- Pouziva sa typicky jedna jedina atomicka konstrukcia/instrukcia – napr compare, swap
- Minimalne prodlevy suvisiace s pristupom k datam
- Neexistuje uviaznutie, je garantovana zivost(program stale bezi, nenastal deadlock, ale to neznamenava ze nehladovie/nestarne)
- Algoritmicke obtiazsie uvazovanie
- Korektnost algoritmu nachlna na optimalizace prekladace

Wait-free procedura

- Procedura, ktorá bez ohľadu na subeh dvoch či viac vlákien, dokončí svoju činnosť v konečnom čase, t.j. neexistuje subeh ktorý by *until* procedure, nekonečne dlho čakal, či vykonávať nekonečne veľa operácií -> nehladovité, nestaré

Lock-free procedura

- Procedúra ktorá garantuje že pri ľubovoľnom subehu veľa súperiach vlákien, vždy ale aspoň jedno vlákno úspešne dokončí svoju činnosť. Niektoré superiace vlákna môžu byť ľubovoľne dlho nutne odkladať dokončenie svojej činnosti

Príklad instrukcií nevhodných k budovaniu lock-free dat. Štruktúra pre viacvlákn. aplikácie

- Test-and-set
- Swap
- fetch-and-add
- Fronty s atomickými operáciami vloženia a výberu

Sú vhodné napríklad pre konštrukciu CAS

*

CAS – compare and swap

- Porovnáva obsah špecifikovanej pamätevej adresy *addr* s očakávanou hodnotou *exp* a v prípade rovnosti nahradí obsah pamätevej adresy novou hodnotou *val*
- O úspechu či neúspechu informuje užívateľa návratovou hodnotou
- Celá procedúra prebehne **atomicky a viditeľne pre ostatné vlákna**

Princíp použitia instrukcie CAS

Postup pri prístupe k zdieľaným dátam

- Precítaj stavajúcu hodnotu zdieľaného objektu
- Pripravím novú hodnotu zdieľaného objektu
- Aplikuj instrukciu CAS

Návratová hodnota

- **True** -> objekt nebol v medzicasе modifikovaný, novo vypočítaná hodnota je platná a uložená v zdieľanom objekte
- **False** -> objekt bol v medzicasе modifikovaný (z iného vlákna), instrukcia CAS nemala žiadny efekt a je nutné celý postup opakovať

Nebezpečie použitia CAS – ABA problem

Kľúčová vlastnosť

- Modifikácia objektu ktorá prebehla medzi načítaním hodnoty objektu a aplikáciou instrukcie CAS nesmie vyprodukovať tú istú hodnotu zdieľaného objektu

Možný chybový scenár

- Hodnota objektu, nacistana vlaknom A za ucelom pouzitia v naslednej nstrukcii CAS, je x
- Pred pouzitim instrukcie CAS vlaknom A, je object modifikovana inymi vlaknami, tj nabyva hodnoty rozne od x
- V okamziku aplikacie instrukcie CAS vlaknom A, ma object opet hodnotu x
- Vlakno A nespозна ze sa hodnota objektu menila
- Nasledna aplikacia instrukcie CAS uspeje

CAS – vykonavanie instrukcii mimo poradie

- Pokial pouzivame CAS na zpristupnenie nejakych dat, je potreba zaistit aby predchadzajuce inicializacie premennych boli uz v okamziku vykonania instrukcie CAS vykonane
- Vyzaduje pouzitie pametovej bariery
- Dotknute premenne musia byt oznacene ako nestale

CAS – cena

- Pouzitie cas odstranilo reziu suvisiacu so zamykanim
- Zustava vsak rezia suvisiaca s coherenciou cache pamete

Programovanie s CAS

Win32

- InterlockedCompareExchange()

Assembler i386(pre x86_64, nutne premenovat edx na rdx)

```
inline int32_t compareAndSwap(volatile int32_t & v, int32_t
exValue, int32_t cmpValue){ asm volatile ("lock; cmpxchg
:%%ecx, %%edx" : "=a" (cmpValue) : "d" (&v), "a" (cmpValue), "c"
(exValue));return cmpValue;}
```

GCC – zaudovane funkcie

- `bool __sync_bool_compare_and_swap (T *ptr, T old, T new, ...)`
- `T __sync_val_compare_and_swap (T *ptr, T old, T new, ...)`

WRRM mapa = priklad

Write Rarely Read Many Mapa

- Zprostredkovava preklad jednej entity na inu(Kluc->hodnota)
- Priklad – kurz koruny k inym menam
 - o Meni sa raz denne
 - o Pouziva sa pri kazdej transakcii

Mozne implementacie s vyuzitim STL

- Map, hash_map

- `assoc_vector`(usporiadane dvojice)

Pouzitie

- `Map<Key,Value> mojeMapa;`

Mutex implementacia *

Implementacia s pouzitim instr. CAS

Operacia citania

- Prebieha zcela bez zamykania

Operacia zapisu

- Vytvorenie kopie stavajúcej mapy
- Modifikacia/pridanie dvojice do vytvorenej kopie
- Atomická zmena novej verzie mapy za predchádzajúcu

Realne omedzenia CAS

- Obecné použitie schémata CAS na `WRRMMap` by vyžadovalo atomickú zmenu relatívne rozsiahlej oblasti pamäte
- HW podpora pre CAS je obmedzená na niekoľko bajtov (typicky jedno, až dva slova procesoru)
- Atomickú zmenu vykonáme cez ukazateľ

*

WRRM map – vlastnosti a problém dealokácie

Prečo je nutná instrukcia CAS a nestaci len `pOld = pNew`

- Interferujú vlákna písary
- Vlákno A urobí kopiu mapy
- vlákno B urobí kopiu mapy, vloží nový kľúč a dokončí operáciu
- vlákno A vloží nový kľúč
- vlákno A nahradí ukazateľ, všetko čo vložilo B je ztratene

Update

- je lock-free, ale není wait-free

Správa pamäte

- ***Update nemože uvoľniť starú kópiu dátovej štruktúry, iné vlákno môže nad dátovou štruktúrou vykonávať operáciu citania***
- Možné riešenie: Garbage collector (JAVA)

Možné riešenie v jazykoch bez garbage collector – odložená dealokácia pamäte

- Miesto delete sa spusti(asynchrone) nove vlakno
- Nove vlakno pocaka 200ms a vykona dealokáciu

Myslienka

- Nove operacie prebiehaju nad novou kopiou, za 200ms sa všetky zapocate operacie nad starou kopiou dokonci a bude bezpecne strukturu dealokovat

Problemy

- Kratkodobe intenzivni prepisovani hodnot alebo vkladanie novych hodnot moze sposobit netrivialne pametove naroky
- Neni garantovane ze sa všetky operacia citania z inych vlakien za dany casovy limit dokoncia

Napad

- Napodobime metodu pouzivanu pri automatickom uvolnovani pameti k tomu, aby sme mohli explicitne dealokova strukturu
- Pocitanie odkazov – s kazdym ukazatelom je zviazane cislo, ktore udava pocet vlakien, ktory tento odkaz este pouzivaju

Modifikacia WRRM mapy

- Procedura *Update* vykona podmienenu dealokáciu, tj, dealokuje object odkazovany ukazatelom, iba pokiaľ žiadne ine vlakno ukazatel nepouziva
- Procedura *Lookup* postupuje tak, ze zvyši citac spojeny s ukazatelom pristupi k structure skrz tento ukazatel, znizi citac po ukonceni prace so strukturu a podmienene dealokuje strukturu

Citac asociovany s ukazatelom $MAP<K,V>^*$

*

- CAS instrukcia nad ukazatelom pDate_
- podmienena dealokacia:

If (pDate_ ->second==0) delete (pDate_ ->first);

Problem v procedure Lookup

- Vlakno A *nacita strukturu* Data(cez *pDate_) a je prerusene
- Vlakno B vlozi kluc, znizi citac a dealokuje *pOld->first
- Vlakno A *zvyši citac*, ale pristupi k neplatnemu ukazatelu

WRRM Map a pocitanie odkazov – CAS2

Problem predchodzej verzie

- Akcia uchopenia ukazatela a zvyšenia odpovedajuceho citaca neboli atomicke

Riesenie

- Pomocou jednej instrukcie CAS je treba prepnut ukazatel a koektnie manipulovat s citacom
- Teoreticky je mozne implementovat CAS pracujuci s viacero strukturami zaroven, samozrejme straca sa efektivita pokiaľ neexistuje HW podpora
- Modern procesory maju podporu pre instrukciu CAS pracujucu s dvoma po sebe ulozenymi slovami procesoru – **CAS2**

*

WRRM Map s vyuzitim CAS2

- Struktura *Data* je tvorena dvoma slovami *ukazatel* a *citac*
- Ukazatel a citac su lozene v pameti vedla seba
- Strukturu je mozne modifikovat pomocou instrukcie CAS2

*

WRRM Map s vyuzitim CAS2 – stale nefunkcne

Otazka

- Dokazeme atomicky realizovat pocitanie odkazov, je teda navrhovane riesenie korektnie?

Problem

- Zvysenie a znizenie citaca procedurou lookup je v zcela nezavislych blokoch. Pokial sa medzi vykonanim tychto blokov realizuje nejaka procedura Update, tak pripocitanie a odcitanie jednotky k citacu prebehne nad inymi ukazatelmi
- Riziko predcasnej dealokacie
- *Ztrata ukazatelov na stare kopie – memory leak*

Riesenie

- Citac spojeny s ukazatelom pouzijeme ako straz
- Procedura update bude vykonavat zmeny struktury len ak ziadne ine vlakno k structure nepristupuje

WRRM Map s vyuzitim CAS2 – realizacia Update

Odkladanie vykonania modifikacie v procedura *Update*

- Atomicke nahradenie ukazatela sa deje v okamihu, kedy su vsetky ostatne vlakna mimo procedure Lookup
- Casove intervaly, po ktorych sa jednotlivé vlakna nachadzaju v procedure Lookup ctenarum sa vsak mozu prekryvat
- Citac po celu dobu existencie ineho vlakna v procedure *Lookup* neklesá na minimalnu hodnotu a procedura *Update* tzv hladovie(Starve)

Optimalizacia procedury *Update*

- pri opakovaných nesuspechoch instrukcie CAS dochádza k opakovanému kopiovaniu pôvodnej štruktúry a následného mazania vytvorenej kopie
- neefektívne opakované kopiovanie lze odstrániť pomocou pomocného ukazateľa *last*

*

WRRM Map – pozorovanie ohľadom realizácie s CAS2

Lookup

- není wait-free, inkrementácia a dekrementácia citaca interferuje s procedúrou update
- volanie procedur update je malo – nevadí

Update

- není wait-free, interferuje s procedúrou Lookup
- volanie procedur lookup je mnoho – problem

Coho sme dosiahli

- WRRM BNTM Mapa
- **Write Rarely Read Many, But Not Too Many**

Dalsie programatorske rozhrania

MCAS

- Rozšírenie standardnej instrukcie CAS pre použitie s ľubovoľne veľkou datovou štruktúrou

Transakčná pamäť

- Pamäť je modifikovaná v jednotlivých transakciách
- Transakcia seskupuje veľa čítaní a zápisu do pamäte – je schopná obsiahnuť komplexnú modifikáciu datových štruktúr
- Základným manipulovateľným objektom je slovo procesoru, tj obsah jednej pamätevej bunky
- Příklad: presun prvku v dynamicky zretazenom zozname

Load-Link(LL)/Store-Conditional(SC)

- Dvojica instrukcia ktorá dohromady realizuje CAS
- LL načíta hodnotu a SC ju prepíše, ale len pokiaľ nebola modifikovaná, pričom za modifikáciu sa považuje aj prepísanie na tú istú hodnotu
- LL/SC stejná sila ako CAS, avšak nemá ABA problém
- HW podpora: Alpha, PowerPC, MIPS, ARM

Problémy

- Zmena kontextu sa v praxi používa za modifikáciu miesta

- Teoreticky není možné realizovat wait-free procedure
- Narocné ladenie

Navrh Lock-free datových struktur

- Je možné navrhnout lockfree datové struktury
- Zaujímavá algoritmika
- Obtiažné pokiaľ chceme deterministicke uvoľňovanie pameti
- Vhodné pre prostredie s garbage collectorom

06 – pokročile rozhrania pre impl. Paral. Aplikácii

Iný spôsob programovania v prostredí so zdieľanou pametou

Nevýhody POSIX threads a lock-free prístupu

- Na prilis nízkej úrovni
- Vhodné pre systémových programátorov
- Prilis zložitý prístup na riešenie jednoduchých vecí

Co by sme chceli

- Paralelná konštrukcia na úrovni programovacieho jazyka
- Prostriedok vhodný pre aplikacných programátorov
- Láhke vyjadrenie bežne používaných paralelných konštrukcií

OpenMP

Myslienka

- Sejmout z programátora bremeno nízkourovňové implementace
- Programátor sa sústreďuje na to, čo chce v programe urobiť, nemá sa zaoberať tým ako sa to stane

Realizace

- Programátor informuje prekladac o zamýšľanej paralelizácii uvedením *znaciek v zdrojovom kóde* a označením blokov
- Pri preklade prekladac sám doplní nízkourovňovú realizáciu paralelizácie

OpenMP ponuka

- Pragma direktívy prekladaca
 - `#pragma omp direktiva [seznam klauzuli]`
- Knihovňové funkcie
- Premenné prostredia

Prekladac kodu

- Prekladac podporujúci standard OpenMP
 - Pri preklade pomocou GCC je nutná voľba `-fopenmp`

- G++ -fopenmp myapp.c
- Podporovane najpouzivanejsimi prekladacmi(aj visual c++)
- Mozno prelozit do sekvencneho kodu → odignoruje znacky pre paralelny kod ak sa nepouzije prepinač

*

Direktiva *parallel*

Použitie

- Strukturovaný blok, tj {...}, nasledujúci sa touto direktívou sa prevedie paralelne
- Mim paralelne bloky sa kód vykonáva sekvencne
- Vlákno ktoré narazi na túto direktívu sa stave hlavným vláknom(master) a má identifikáciu vlákna rovnú 0

Podmienene spustenie

- Klausula if(vyraz typu bool)
- Ak sa vyhodnotí výraz na false direktiva *parallel* sa ignoruje a nasledujúci blok je vykonaný iba v jednej kopii

Stupen paralelizmu

- Počet vlákien
- Prednastavený počet špecifikovaný premennou prostredia
- Klausula *num_thread*(výraz typu int)

Direktiva *parallel* – datová lokalita

Klausula *private* (seznam premennych)

- Vymenované premenné sa zduplikujú a stanú sa lokálnymi premennými v každom vlákne

Klausule *firstprivate* (seznam premennych)

- Viz *private*, s tým, že všetky kopie premenných sú inicializované na hodnotu originalnej kopie

Klausule *shared* (zoznam premennych)

- Vymenované premenné budú explicitne existovať v iba jednej kopii
- Prístup k zdieľaným premenným nutne serializovať

Klausule *default* ([*shared*/*none*])

- *Shared*: všetky premenné sú zdieľané, pokiaľ není uvedené inak
- *None*: vynucuje explicitné uvedenie každej premennej v klauzuli *private* alebo v klauzuli *shared*

Direktiva *parallel* – redukce

Klauzule ***reduction*** (*operator: seznam premennych*)

- Pri ukončení paralelného bloku sú vymenované private premenne zkombinované pomocou uvedeného operatoru
- Kopie uvedených premenných, ktoré sú platné po ukončení paralelného bloku, sú naplnené výslednou hodnotou
- Premenne musia byť skalárneho typu (nesmie byť pole, štruktúra, etc.)
- Použiteľné operatory: +, *, &, |, ^, && a ||

Direktiva ***for***

Použitie

- Iterácia nasledujúceho for cyklu budú prevedené paralelne
- Musí byť použité v rámci bloku za direktívou *parallel* (inak prebehne sekvencne)
- Možný zkratený zápis #pragma omp parallel for

Klauzule ***private, firstprivate, reduction***

- Stejne ako pre direktívu parallel

Klauzule ***lastprivate***

- Hodnota privatej premennej vo vlákne spracovávajúcej poslednú iteráciu for cyklu je uložená do kopie premennej platnej po skončení cyklu

Klauzule ***ordered***

- Bloky označené direktívou ordered v tele paralelne provedeného cyklu sú prevedené v tom poradí, v akom by boli prevedené sekvencným programom
- Klausule ordered je povinná, pokiaľ telo cyklu obsahuje ordered bloky

Klauzule ***nowait***

- Jednotlivé vlákna sa nesynchronizujú po vykonaní cyklu

Klauzule ***schedule*** (*typ planovani* [, *velkost*])

- Urcuje ako budú iterácie rozdelené/mapované medzi vlákna
- Implicitné plánovanie je závislé na implementácii

Direktiva ***for*** – plánovanie iterácii

Static

- Iterácie cyklu rozdelené do blokov o špecifickej veľkosti
- Bloky staticky namapované na vlákna (round robin)

- Round robin – priradi procesu/vlaknu kvantum casu, po ktore moze byt process/vlakno na procesore. Po ubehnuti casu je spustene ine vlakno/proces
- Pokial neni uvedena velkost, iteracie rozdelené medzi vlakna rovnomerne(ak je to mozne)

Dynamic

- Bloky iteracii cyklov v pocte specifikovanim parametrom *velkost* pridelovane vlaknam na ziadost, tj v okamihu kedy vlakno dokoncilu svoju predchadzajucu pracu
- Defaultna velkost bloku je 1

Guided

- Bloky iteracii maju velkost proportionalnu k poctu nezpracovanych iteracii vydelenym poctom vlakien
- Specifikovana velkost *k* udava minimalnu velkost bloku(defaultne 1)
- *Priklad:*
 - $K=7$, 200 volnych iteracii, 8 vlakien
 - Velkost bloku: $200/8=25$, $175/8=21$, ..., $63/8=7$, ...

Runtime

- Typ planovania urceny az za behu premennou *OMP_SCHEDULE*

Direktiva *sections*

Pouzitie:

- Strukturovane bloky, kazdy oznaceny direktivou *section*, mozu byt vramci bloku oznacenym direktivou *sections* vykonavane paralelne
- Mozny zkrateny zapis `#pragma omp parallel sections`
- Umoznuje definovat rozny kod pre rozna vlakna

Klauzule *private*, *firstprivate*, *reduction*, *nowait*

- Rovnake ako v predoslych pripadoch

Klauzule *lastprivate*

- Hodnoty privatnych premennych v poslednej sekcii(podla zapisu kodu) budu platne po skoncení bloku *sections*

*

Vnoroenie direktiv *parallel*

Nevnorený paralelizmus

- Direktiva *parallel* urcuje vznik oblasti paralelného vykonavania

- Direktivy for a sections urcuje ako bude praca mapovana na vlakna podla rodicovskej direktivy *parallel*

Vnorený paralelizmus

- Pri nutnosti paralelizmu v rámci paralelného bloku, je treba znovu uviesť direktivu *parallel*
- Vnorovanie je podmienené nastavením premennej prostredia OMP_NESTED(hodnoty TRUE, FALSE)
- Typické použité vnorené for-cykly
- Obecné vnorovanie direktív v OpenMP pomerne komplikované
- Može vzniknúť viac vlákien než CPU jadier – systém príliš zatažený tým, že strieda vlákna na jadrách a nie je teda vhodný vnorený paralelizmus vždy používať

Direktiva *barrier*

Bariera

- Miesto, ktoré je dovolené prekročiť až k nemu dorazia všetky ostatné vlákna
- Direktiva bez klauzuly, tj `#pragma omp barrier`
- Vztahuje sa k štrukturalne **najbližšej** directive *parallel*
- Musí byť volané **všetkými** vláknami v odpovedajúcom bloku direktivy *parallel*
 - Nefunguje teda zavolať na napríklad 2 zo 4
 - Da sa obísť tým, že sa v *parallel* bloku vytvorí blok na 2 vlákna, kde už to bude fungovať

Poznámka ku kodovaniu

- Direktivy prekladače nie sú súčasťou jazyka
- Je možné, že v rámci prekladu bude vyhodnotený blok, v ktorom je umiestnená direktiva *barrier*, ako nevykonateľný blok a odpovedajúci kód nebude vo výslednom spustiteľnom súbore vôbec prítomný
- Direktivu *barrier* je nutné umiestniť v bloku, ktorý sa bezpodmienečne vykoná (zodpovednosť programátora)

Direktiva *single a master*

Direktiva *single*

- V kontexte paralelne vykonávaného bloku je nasledujúci štrukturalný blok vykonaný iba jedným vláknom, pričom není určené ktorým

Klauzule *private, firstprivate*

Klauzule *nowait*

- Pokiaľ není uvedená, tak na konci štrukturalného bloku označeného direktívou *single* je vykonaná bariera

Direktiva *master*

- Specialny pripad direktivy single
- Tym vlaknom, ktore vykona strazeny blok, bude hlavne(master) vlakno

Direktiva *critical a atomic*

Direktiva *critical*

- Nasledujuci strukturovany blok je chapany ako kriticka sekcia a moze byt vykonany maximalne jednym vlaknom v danom case
- Kriticka sekcia moze byt pomenovana, subezne je mozne vykonavat kod v kritickyh sekciach s inym nazvom
- Pokial neni uvedene inak, pouzije sa implicitne meno
- `#pragma omp critical [(name)]`

Direktiva *atomic*

- Nahradzuje kriticku sekciu nad jednoduchymi modifikaciami(update) premennych v pameti
- Atomicita sa aplikuje na jeden nasledujuci vyraz
- Obecne vyraz musi byt jednoduchy(len *load a store*)
- Neatomizovatelny vyraz: $x = y = 0$ -> dve ulozenia do pameti

Direktiva *flush*

Problem(nestale premenne)

- Modifikacie zdielanych premennych v jednom vlakne moze ostat skryta ostatnym vlaknam

Riesenie

- Explicitna direktiva pre kopirovanie hodnoty premennej z registru do pameti a zpet
- `#pragma omp flush [(seznam)]`

Pouzitie

- Po zapise do zdielanej premennej
- Pred citanim obsahu zdielanej premennej
- Implicitne v miestach bariery a konca blokov(pokial niesu bloky v rezime nowait)

Direktiva *threadprivate a copyin*

Problem(thread-private data)

- Pri statickom mapovani na vlakna je drahe pri opakovanom vznku a zaniku vlakien vytvarat kopie privatnych premennych
- Obcas chceme private globalne premenne

Riesenie

- Perzistentne private premenne – preziju zanik vlakna
- Pri znovuvytvorení vlakna sa premenne znovupouziju
- `#pragma omp threadprivate` (seznam)

Obmedzenia

- Nesmie sa použiť dynamicke planovanie vlákien
- Počet vlákien v paralelných blokoch musí byť zhodný

Direktiva *copyin*

- Ako `threadprivate`, ale s inicializáciou
- Viz `private` versus `firstprivate`

OpenMP knihovnové funkcie – počet vlákien

```
void omp_set_num_threads (int num_threads)
```

- Specifikuje koľko vlákien sa vytvorí pri ďalšej directive `parallel`
- Musí byť použité pred samotnou konštrukciou `parallel`
- Je prebité klauzulou `num_threads`, pokiaľ je prítomná
- Musí byť povolené dynamicke modifikovanie procesov (`OMP_DYNAMIC`, `omp_set_dynamic()`)

```
int omp_get_num_threads ()
```

- Vracia počet vlákien v tíme štruktúrne najbližšej direktívy `parallel`, pokiaľ neexistuje, vracia 1

OpenMP knihovnové funkcie – počet vlákien a procesorov

```
int omp_get_max_threads ()
```

- Vracia maximálny počet vlákien v tíme.

```
int omp_get_thread_num ()
```

- Vracia unikátny identifikátor vlákna v rámci tímu.

```
int omp_get_num_procs ()
```

- Vracia počet dostupných procesorů, ktoré môžu v danom okamžiku participovať na vykonávaní paralelného kódu.

```
int omp_in_parallel ()
```

- Vracia nenula pokiaľ je voláno v rozsahu paralelného bloku.

OpenMP knihovnové funkcie – kontrola vytvárania vlákien

```
void omp_set_dynamic (int dynamic_threads)
```

```
int omp_get_dynamic()
```

- Nastavuje a vracia, či je programátorovi umožnené dynamicky meniť počet vlákien vytvorených pri dosiahnutí direktívy `parallel`
- Nenulová hodnota `dynamic_threads` znamená **povoľeno**

```
void omp_set_nested (int nested)
```

```
int omp_get_nested()
```

- Nastavuje a vracia, či je povolený vnorený paralelizmus
- Pokiaľ nie je povolené, vnorené paralelne bloky sú serializované

OpenMP knihovnové funkcie – mutexy

```
void omp_init_lock (omp_lock_t *lock)
```

```
void omp_destroy_lock (omp_lock_t *lock)
```

```
void omp_set_lock (omp_lock_t *lock)
```

```
void omp_unset_lock (omp_lock_t *lock)
```

```
int omp_test_lock (omp_lock_t *lock)
```

```
void omp_init_nest_lock (omp_nest_lock_t *lock)
```

```
void omp_destroy_nest_lock (omp_nest_lock_t *lock)
```

```
void omp_set_nest_lock (omp_nest_lock_t *lock)
```

```
void omp_unset_nest_lock (omp_nest_lock_t *lock)
```

```
int omp_test_nest_lock (omp_nest_lock_t *lock)
```

- Inicializuje, ničí, blokuje a čaká, odemyká a testuje
- – normálny a rekurzívny mutex.

Premenné prostredia

OMP_NUM_THREADS

- Špecifikuje defaultný počet vlákien, ktoré sa vytvoria pri použití direktívy `parallel`

OMP_DYNAMIC

- Hodnota TRUE, umožňuje za behu dynamicky meniť počet vlákien

OMP_NESTED

- Povoluje hodnotou TRUE vnorený paralelizmus
- Hodnotou FALSE špecifikuje že vnorené paralelné konštrukcie budú serializované

OMP_SCHEDULE

- Udať defaultné nastavenie mapovania iterácií cyklu na vlákna
- Príklady hodnôt: "static, 4", dynamic, guided

Intel's Thread Building Blocks (TBB)

Co je Intel TBB

- TBB je c++ knižnica pre vytváranie viacvláknových aplikácií
- Založená na princípu zvanom *generic programming*
- Vyvinuté synergickým spojením Pragma direktív(OpenMP), štandardnej knižnice šablón(STL, STAPL) a programovacích jazykov podporujúcich prácu s vláknami(Threaded-C, Cilk)

Generic Programming

- Vytváranie aplikácií špecializáciou existujúcich predpripravených obecných konštrukcií, objektov a algoritmov
- Da sa najst' v objektovo orientovaných jazykoch – C++, JAVA
- V c++ sú obecnou konštrukciou šablóny(Templates)
 - o Queue<int>
 - o Queue<Queue<Char>>

Použitie TBB

Vlastnosti Intel TBB

- Knižnica, implementovaná s využitím štandardného c++
- Nepožaduje podporu špeciálneho jazyka či prekladáča
- Podporuje vnorený paralelizmus, potažmo je možné stavať zložitejšie paralelné systémy z menších paralelných komponentov
- Cieľom použitia je nechať programátora špecifikovať úlohy k paralelnému vykonaniu, nie nútiť ho popisovať čo a ako robia jednotlivé vlákna

Možnosti TBB

TBB poskytuje šablóny pre

- Paralelizáciou iterácií jednotlivých cyklov – dátový paralelizmus
- Definíciu vlastných paralelne prístupovaných dátových štruktúr
- Využitie nízkourovnových HW primitív
- Zamykanie prístupov do kritickej sekcie v rôznych podobách

- Lahku definiciu paralelnych subeznych uloh
- Skalovatelnu alokaciju pameti

Princip datovej paralelizacie v TBB

Paralelny for-cyklus

- Je dana mnozina nezávislych indexov, tzv rozsah(range)
- Pre kazdy index z mnoziny je vykonane telo cyklu

Paralelny for-cyklus v TBB

- Sablona, ktora ma dva parametry – rozsah a telo cyklu
- Sablona zaisti vykonanie tela cyklu pre vsetky indexy v specifikovanom rozsahu
- Rozsah je deleny na pod-rozsahy. Paralelizmu dosiahnute vykonavaním tela cyklu nad jednotlivymi pod-rozsahmi

*

Koncept delenia

- Instancie niektorých tried je nutné za behu (rekurzívne) delit
- Zavádza sa nový typ konštruktoru, deliaci konštruktor:
 - `X::X(X& x, split)`
- Deliaci konštruktor rozdelí instanciu triedy X na dve časti, ktoré dohromady dávajú pôvodný objekt. Jedna časť je priradená do x, druhá časť je priradená do novo vzniknutej instance
- Schopnosť deliť sa musia mať najmä rozsahy, ale taktiež triedy ktorých instance bežia paralelne a pritom nejakým spôsobom interagujú, napr triedy realizujúce paralelnú redukciu

Split

- Špeciálna trieda definovaná za účelom odlíšenia deliaceho konštruktoru od kopirovacieho konštruktoru

Koncept rozsahu

Požiadavky na triedu realizujúcu rozsah

- Kopirovací konštruktor
 - `R::R (const R&)`
- Deliaci konštruktor
 - `R::R (constR&, split)`
- Destruktor
 - `R::~~R ()`
- Test na prázdnotu rozsahu
 - `Bool R::empty() const`
- Test na schopnosť ďalšieho rozdelenia

- Bool R::is_divisible() const

Preddefinovane ssablony rozsahov

- Jednodimenzionalne: blocked_range
- Dvojdimonzionalne: blocked_range2d

TBB: blocked range

Blocked_range

- Template<typename Value> class blocked_range;
- Reprezentuje nadalej delitelny otvoreny interval [l,j)

Poziadavky na triedu *Value* specializujucu *blocked_range*

- kopírovací konstruktor
 - Value::Value (const Value&)
- Destruktor
 - Value::~~Value ()
- Operátor porovnání
 - bool Value::operator<(const Value& i, const Value& j)
- Počet objektů v daném rozsahu (operátor-)
 - size_t Value::operator-(const Value& i, const Value& j)
- k-tý následný objekt po i(operátor+)
 - Value Value::operator+(const Value& i)

Použitie blocked_Range<Value>

- Najdoležitejšou metódou je konstruktor
- Konstruktor špecifikuje interval rozsahu a veľkosť najväčšieho ďalej nedeliteľného sub intervalu
- Blocked_range(Value begin, Value end, [size_t gransize])

Typická specializácia

- Blocked_range <int>
- Příklad blocked_range<int>(5,17,2)
- Příklad blocked_range<int>(0,11)

TBB: parallel_for

Parallel_for<Range, Body>

- template<typename Range, typename Body>
- void parallel_for(const Range& range, const Body& body);

Poziadavky na triedu realizujuce telo cyklu

- konstruktor
- destructor

- aplikator tela cyklu na dany rozsah – operator()

TBB: parallel_reduce

Parallel_Reduce<range,body>

Poziadavky na triedu realizujuce telo redukcie

- deliaci konstruktor
- destructor
- funkcia realizujuca redukciu nad danym rozsahom – operator()
- funkcia realizujuca redukciu hodnot z roznych rozsahov

Moznosi delenia

Trieda *Partitioner*

- paralelne konstrukcie maju tretí voliteľný parameter, ktorý špecifikuje stratégiu delenia rozsahu
- parallel_for<range,body,partitioner>

Preddefinované stratégie

- simple_partitioner
 - o rekurzívne deli rozsah až na ďalej nedeliteľné intervaly
 - o pri použití blocked_range je voľba grainsize kľúčová pre vyváženie potenciálu a rezie paralelizácie
- auto_partitioner
 - o automatické delenie, ktoré zohľadňuje zataženie vlákien
 - o pri použití blocked_range volí rozsahy väčšie než je grainsize a tieto deli iba do tej doby, než je dosiahnuté rozumné vyváženie zataženia. Voľba minimálnej veľkosti grainsize nespôsobí nadbytočnú reziu spojenú s paralelizáciou

Paralelne prístupované kontajnery – *vector* a *queue*

Concurrent_queue

- template<typename T> concurrent_queue
- Fronta, ku ktorej môže súbežne prístupovať viacero vlákien
- Veľkosť fronty je daná počtom operácií vloženia bez počtu operácií vyberu, záporná hodnota znamená čakajúce operácie vyberu
- Definuje sekvencné iterátory, neodporúča sa ich používať

Concurrent_vector

- Template<typename T>concurrent_vector

- Zvacsovateľne pole prvkov, ku ktorým je možné subezne pristupovať z viacerých vlákien a vykonávať subezne zvacsovania polí a pristup k už uloženým prvkom
- Nad vektorom sa dá definovať rozsah a vykonávať skrz neho paralelne operácie s prvkami uloženými v poli

Paralelne pristupovane kontajnery – hash_map

Concurrent_hash_map

- Mapa, v ktorej je možné paralelne hľadať, mazat a vkladať

Požiadavky na triedu HashCompare

- Kopirovací konštruktor
- Destructor
- Test na ekvivalenciu objektov
- Výpočet hodnoty hashovacej funkcie

Objekty pre pristup k datam v *concurrent_hash_map*

- Prístup k parom kľuč-hodnota je skrz pristupovacie triedy
- *Accessor* – pre prístup v režime read/write
- *Const_accessor* – pre prístup iba v režime read
- Použitie pristupovacích objektov umožňuje korektný paralelný prístup k zdieľaným dátam

*

Metody pre prácu s *concurrent_Hash_map*

- Find – dva typy -> pre konštantný accessor a normálny accessor (da sa modifikovať)
- Insert
- erase

Dalsie sposoby použitia

- Iteratory pre prechádzanie mapy
- Dajú sa definovať rozsahy a s nimi pracovať paralelne

C++11

C++11 a vlakovanie

Pozorovanie

- C++11 má definované príkazy pre podporu vlákien
- Není treba použiť externé knižkovne ako je POSIX Thread

Ako je to možné

- C++11 definuje virtualne viacjadrový výpočetný stroj a teda ráta sa s tým že jazyk bude bežať na viacjadrovom výpočetnom stroji
- Všetka semantika príkazov sa odkazuje na tento virtuálny výpočetný stroj
- Príkazy s podporou vlákien môžu byť súčasťou jazyka
- Prenos semantiky z virtuálneho výpočetného stroja na reálny HW je na zodpovednosť prekladáča

Vlákna a mutexy v C++11 *

Zamykanie v C++11

Potencionálne riziko uviaznutia

- Jazyk s plnou podporou mechanizmu výnimiek
- Vyhodenie výnimky v okamihu keď je vlákno v kritickej sekcii (uvnit' mutexu) pravdepodobne spôsobí, že nebude vláknom volaná metóda odomykajúca zámek spojený s kritickej sekcii

Reseni

- Využitie princípu RAII (**Resource acquisition is initialization**) a OOP
- Zamčenie mutexu realizované vytvorením lokálnej inšancie hodnej preddefinovanej zamykacej triedy
- Odomknutie umiestnené do deštruktora tejto triedy
- Deštruktor je vykonaný v okamihu opustenia rozsahu platnosti daného objektu

RAII Zamykanie v C++11

Trieda *lock_guard*

- Obalenie štandardného zámku v RAII style
- Mutex na pozadí nejde predať inému vláknu, nevhodné pre podmienkové premenne
- Príklad použitia *

Trieda *unique_lock*

- Obecnejší predateľný RAII obalenie mutexu
- Doporučené pre použitie s podmienkovými premennými

Podporované aspekty

Podpora vlákna v C++11

- Vlákna
- Mutexy a RAII zámky
- Podmienkové premenne
- Zdieľané futures (miesta kde sa uložia dosiaľ nespocítané hodnoty)

Atomicita zápisu

Neatomicky

Int x,y;

Thread 1

X=17

Y=37

Thread 2

Cout << y << " ";

Cout << x << endl;

Nema definovane chovanie

Spravne atomicky

Atomic<Int> x,y;

Thread 1

X.store(17)

y.store(37)

Thread 2

Cout << y.load() << " ";

Cout << x .load()<< endl;

Chovanie je definovane, moze vystupy 0 0, 0 17, 37 17

Pametovy model v c++11

- Implicitne chovanie zachovava sekvenčnu konzistenciu(automaticky vklada odpovedajuce pametove bariery)
- Riziko neefektivneho kodu

Priklad 1

```
atomic<int> x, y;
```

Thread 1

```
x.store(17,memory_order_relaxed);
```

```
y.store(37,memory_order_relaxed);
```

Thread 2

```
cout << y.load(memory_order_relaxed) << " ";
```

```
cout << x.load(memory_order_relaxed) << endl;
```

Semantika povoľuje v tomto prípade I výstup 37 0

Priklad 2

```
atomic<int> x, y;
```

Thread 1

```
x.store(17,memory_order_release);
```

```
y.store(37,memory_order_release);
```

Thread 2

```
cout << y.load(memory_order_acquire) << " ";
```

```
cout << x.load(memory_order_acquire) << endl;
```

Acquire nepreusporiada operace load, Release-store

Ine pristupy

Paralelny for cyklus

- Najcastejsi a najjednoduchsia metoda paralelizacie
- Datova paralelizacia

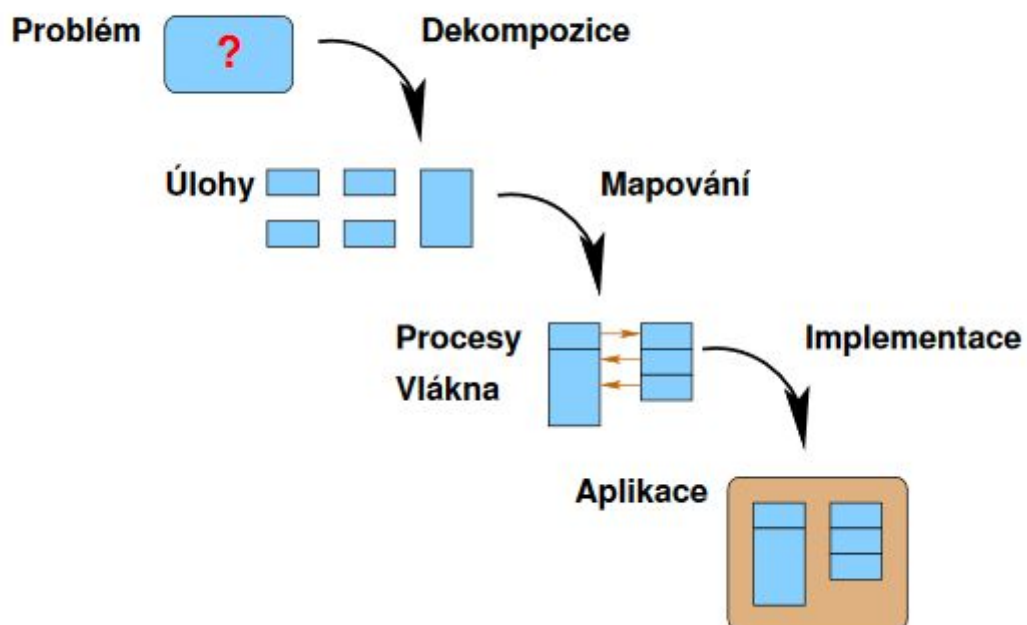
07 – principy navrhnu paralelnych algoritmov

Vice-prace programatora paral. Aplikacii

- identifikovat subezne vykonatelne cinnosti a ich zavislosti
- mapovat subezne proveditelne casti prace do procesu
- zaistit distribuciu vstupnych, vnutornych a vystupnych dat
- spravovat subezny pristup k datam a zdielanym prostriedkom
- synchronizovat jednotlivé procesy v roznych stadiach vypoctu tak, ako vyzaduje paralelny algoritmus
- mat znalost pridavnych programatorskych prostriedkov suvisiacich s vyvojom paralelnych algoritmov

Zaklady navrchu paralelych algoritmov

Navrh a realizace paralelného systému



Dekompozice a ulohy

Dekompozice

- process rozdelenia celej vypocetnej ulohy na podulohy
- niektore podulohy mozu byt vykonavane paralelne

(pod)ulohy

- jednotky vypoctu ziskane dekompoziciou
- po vycleneni sa povazuju za dalej nedeliteľne
- maju uniformnu/neuniformnu velkost
- su definovane v dobe kompilace/za behu programu

Příklad

- násobení matice $A(n \times n)$ vektorem B

Zavislosti uloh

Graf zavislosti

- zachycuje zavislosti vykonavanych uloh
- definuje relativne poradie vykonavanych uloh (ciastocne usporiadanie)

Vlastnosti a vyuzitie grafu

- orientovany acyklicky graf
- graf moze byt nespojity ci dokonca prazdny
- uloha je pripravena k spusteniu, pokiaľ ulohy na ktorých zavisí, dokončili svoj výpočet (topologické usporiadanie)

Priklady zavislosti

- poradie obliekania svrsku
- paralelní vyhodnocování výrazu
 - $v \text{ AND } x \text{ AND } (y \text{ OR } z)$

Granularita a stupen subeznosti

Granularita

- počet uloh na ktoré sa problém dekomponuje
- mnoho malých uloh – jemnozrnná granularita (fine-grained)
- málo väčších uloh – hrubozrnná granularita (coarse-grained)
- každý problém má vnútornú hranicu granularity

Stupen subeznosti

- maximálny počet uloh, ktoré môžu byť vykonávané súčasne
- limitom je vnútorná hranica granularity

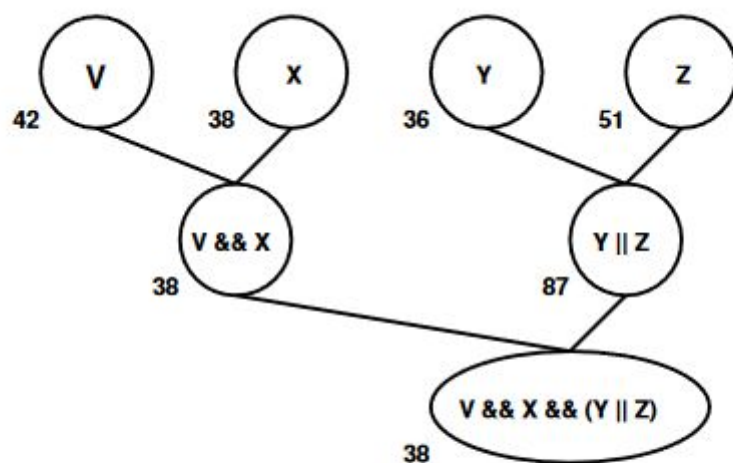
Priemerný stupen subeznosti

- závislý na grafe závislosti a granularite
- má najmä množstvo práce asociované k uzlom grafu
- **kritická cesta** – cesta, na ktorej je súčet prác maximálny
- **priemerný stupen subeznosti** - je podiel celkového množstva práce v celkovom množstve prác na kritickej ceste
- udáva maximálne zrýchlenie, pokiaľ je cieľová platforma schopná vykonávať subezne maximálny stupen subeznosti uloh

Pozorovania

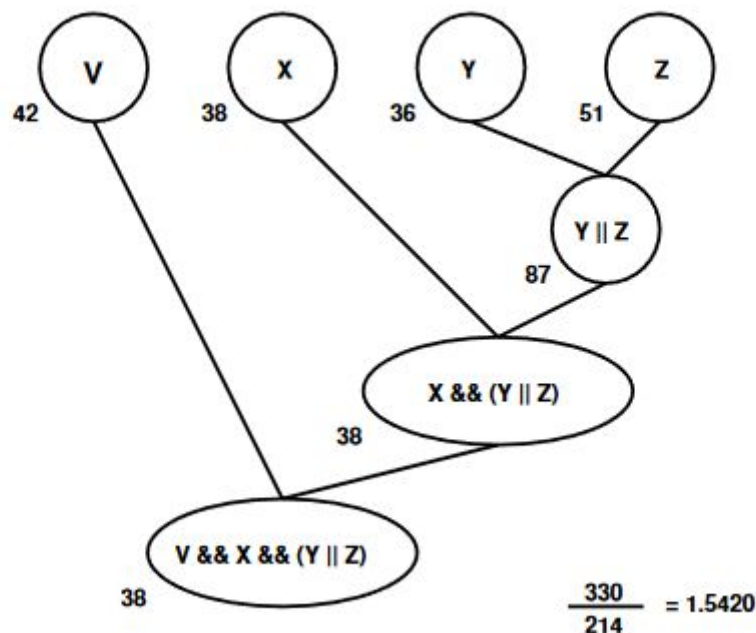
- zjemňovanie dekompozície môže zvýšiť stupen subeznosti
- čím menej práce je na kritickej ceste tým väčší je potenciál pre paralelizáciu

Priemerny stupen subeznosti – Varianta 1



$$\frac{330}{176} = 1.8750$$

Priemerny stupen subeznosti – varianta 2



Interakce uloh – dalsi omezeni paralelizace

Interakce uloh

- nezávislé úlohy mohou vzájemně komunikovat
- obojsměrná komunikace může snížit stupeň subeznosti (úlohy musí co-existovat v stejný okamžik)
- komunikace uloh – **neorientovaný graf interakcí**
- graf interakcí pokrývá graf závislosti (overení splnění předpokladu pro spuštění ulohy je forma interakce)

Příklad jednosměrné komunikace

- násobení matice vektorem ($y = Ab$)
- dekompozice na nezávislé úlohy podle řádku matice A
- prvky vektoru b se čítají ze všech úloh, je nutné je vhodně distribuovat k jednotlivým úlohám

Techniky dekompozice

Dekompozice

- fundamentální technika v návrhu paralelních algoritmů

Klasifikace technik dekompozice

- obecné techniky dekompozice
 - o rekurzivní
 - o datová
- specializované techniky dekompozice
 - o průzkumná
 - o spekulativní

Hybridní přístup

- kombinace výše uvedených technik

Rekurzivní dekompozice

Vhodná pro problém typu rozděli a panuj

Princip

- problém se dekomponuje na podúlohy tak, aby se jednotlivé úlohy mohly mohly dekomponovat rovnakým způsobem jako rodičovská úloha
- někdy je třeba restrukturalizovat úlohu

Příklad

- Quicksort
 - o Průběh se volba pivotu
 - o Rozdělení pole na prvky menší než a větší rovně než
 - o Rekurzivně se opakuje dokud je množina prvků neprázdná
- Hledání minima v lineárním poli
 - o Princip pulení prehladavaneho pole
 - o Typický příklad restrukturalizace výpočtu

Datová dekompozice

Základní princip

- Data se rozdělí na části (data partitioning)
- Úlohy se vykonávají subežně nad jednotlivými částmi dat

Datová dekompozice podle místa

- Vstupní data
- Výstupní data
- Vnitřní data
- Kombinace

Mapování dat na úlohy

- Funkce identifikuje vlákno zodpovědné za zpracování dat

Úlohy typu "embarrassingly parallel"

- Trivialna datova dekompozícia na dostatečný počet zcelá nezávislých, vzajomne nekomunikujúcich úloh

Priezkumova dekompozícia

Princíp

- Specializovaná technika paralelizácie
- Vhodná pre prehľadávajúce úlohy
- Prehľadávaný priestor sa rozdelí podľa smeru hľadania

Vlastnosti

- Pri znalosti prehľadávaného stavového priestoru sa dá dosiahnuť optimálneho vyváženia a zatazenia procesoru
- Na rozdiel od datovej dekompozície, úloha končí okamžite keď je najdené požadované
- ***Množstvo vykonanej práce sa líši od sekvencnej verzie***
- V prípade že graf není strom, je treba riešiť problém s opakujúcimi konfiguráciami (riziko nekonečného výpočtu)

Príklad

- Riešenie hlavolamu “patnact”

Spekulatívna dekompozícia

Princíp

- Specializovaná technika paralelizácie
- Vhodná pre úlohy so sekvenciou datovo závislých podúloh
- Úloha ktorá čaká na výstup predchádzajúcej úlohy sa spustí nad všetkými možnými výstupmi (výstupy predchádzajúcej úlohy)

Vlastnosti

- Vykonáva sa zbytočná práca
- Nemusi byť vo výsledku rýchlejšia ako serializovaná verzia
- Vhodná pre úlohy, kde určitá hodnota medzivýsledku má veľkú pravdepodobnosť
- Vznika potenciálny problém pri prístupe k zdrojom (niektoré zdroje nemusia byť zdieľané v prípade sekvencného vykonávania úloh)

Priklady

- Spekulatívne vykonávanie kódu – vetvenie
- Odhadovanie výsledkov operácií v CPU, zahodenie nesprávneho

Zriedka používané

Hybridná dekompozícia

Priklad – hladanie minima v poli

- Sekvenca verzia najde minimum v $O(n)$
- Pri pouziti datovej a rekurzivnej dekompozicie sa da trvanie tejto ulohy skratit na $O(n/p + \log(p))$
- Vstupne pole sa datovo dekomponuje na p rovnakych casti
- Najdu sa minima jednotlivych casti v case $O(n/p)$
- Vysledky sa zkombinuju v case $O(\log(p))$
- Teoreticky sa da pri dostatočnom počte procesorov najst minimum v case $O(\log(n))$

Pozorovanie

- Postup pouzity vo vyššie uvedenom priklade sleduje obecne schema, ktora je oznacovane ako **MAP-REDUCE**

Techniky mapovania a vyrovnavania zateze

Mapovanie uloh na vlakna/procesy

Mapovanie

- Priradovanie uloh jednotlivym vlaknam/procesom
- Optimalne mapovanie berie v ohľad grafy závislosti a interakcie
- Ovplyvnuje výkon aplikacie
- Naivne mapovanie (uloha=process/vlakno)
- Je to Mentalny process, tak ako dekompozicie

Cile mapovania

- Minimalizovat celkovy čas riesenia celej ulohy
 - o redukovať omeškanie spôsobene čakanim(idling)
 - o redukovať zataz spôsobenu interakciou
 - o redukovať reziu spustania, ukončovania a prepínania
 - o vyrovnavat zataz na jednotlivé procesy
- maximalizovat subeznosť
- minimalizovat zatazenia system(zatazenie datovych ciest)
- vyuzivat dostupnosť zdrojov pouzitych predchadzajucou ulohou

Charakteristiky uloh, ktore ovplyvnuju mapovanie

Sposob zadania ulohy

- **staticke zadanie uloh**
 - o dekompozicia problem na ulohy je dana v dobe kompilacie, pripadne je priamo odvodená od vstupnych dat
- **dynamicke zadanie uloh**

- nove ulohy su vytvarane za behu aplikacie podla priebehu vypoctu pripade ako dosledok vykonavanie povodne zadanych uloh

velkost ulohy

- relativne mnozstvo casu potrebne k dokonceniu ulohy
- uniformna vs neuniformna
- dopredna znalost/neznalost

velkost dat asociovaných k ulože

- snaha o zachovanie locality dat
- rozne data maju roznu rolu a velkost (vstupne/vystupne data u hlavolamu patnact)

Staticke vs dynamicke

- staticke: prebiehaju v preddefinovanom casovom intervale, medzi predom znamou mnozinou uloh
- dynamicke: pokial predom nepozname pocet interakcii, casovy ramec interakcii alebo participujuce ulohy

Dalsie charakteristiky

- jednosmerna versus obojsmerna interakcia
- mod pristupu k datam -> read-Only vs read-write
- pravidelne versus nahodile interakce

Rezie suvisiaca s mapovaním do roznych vlakien/procesov

- usposobenie aplikacie pre neocakavanu interakciu
- pripravenost dat k odoslaniu/ adresata k priatiu
- riadenie pristupu k zdielanym zdrojom
- optimalizacia aplikacie pre redukciu prodlev

Schemata pre staticke mapovanie

Mapovanie zalozene na rozdeleni dat

- blokova distribucia
- cyklicka a blokovo-cyklicka distribucia
- nahodna distribucia blokov
- delenie grafu

Mapovanie zalozene na rozdeleni uloh

- delenie podla grafu zavislosti uloh
- hierarchicke delenie

Mapovanie zalozene na rozdeleni dat

Blokova distribucia datovych poli

- procesy spojené s datami rozdelenými na súvislé bloky
- bloky môžu byť viacrozmerné (redukcia interakcií)
- príklad
 - o násobenie matic $A \times B = C$
 - o delenie matic C na 1- a 2-rozmerné bloky

Cyklická a blokovo cyklická distribúcia dátových poli

- nerovnomerné množstvo práce spojené s jednotlivými prvkami → blokovo delenie spôsobuje nerovnomerné zaťaženie
- blokovo-cyklická distribúcia: delenie na menšie diely a cyklické priradenie procesom (round robin)
- zmenšovanie bloku vedie k cyklickej distribúcii (blok je atomický prvok dátového pole)

Náhodná distribúcia bloku

- zaťaž súvisiaca s prvkami pola vytvára pravidelné vzory
- → spätná distribúcia v cyklickom rozdelení
- Náhodné priradenie bloku procesom

Grafové delenie

- Pre prípady, keď je nevhodné organizovať dáta do poli (napríklad drátové modely 3D objektov)
- Dáta organizované ako graf
- Optimálne delenie
 - o Stejný počet vrcholov v jednotlivých častiach
 - o Čo možno najmenší počet hrán medzi jednotlivými časťami → čo najmenej komunikácie medzi vláknami
 - o NP úplný problém

Mapovanie založené na rozdelení úloh

Princíp

- Graf závislosti úloh
- Grafové delenie (NP-úplné)

Špeciálne prípady pre konkrétny tvar grafov

- Binárny strom (rekurzívna dekompozícia)

Hierarchické mapovanie

- Úlohové delenie neberie v úvah neuniformitu úloh
- Shlukovanie úloh do nad-úloh
- Definuje hierarchie (vrstvy)
- Iné mapovacie a dekompozičné techniky na jednotlivých vrstvách

Schemata pre dynamicke mapovanie

Motivacia

- Staticke mapovanie nedostatočne, pretože charakteristiky uloh nie sú známe v dobe prekladu

Centralizovane schemata dynamickeho mapovania

- Ulohy sú zhromažďované v jednom mieste
- Dedikovaná uloha pre priradovanie uloh procesom
- Samo-planovanie
 - o Akonáhle proces dokončí ulohu vezme si ďalšiu
- Blokove planovanie
 - o Prístup k zhromažďovaniu uloh môže byť úzkym miestom
 - o -> pridelovanie uloh po blokoch

Priklad

- Triedenie prvkov v $n \times n$ matici A
- `For(i=1; i<n; i++) newtask(sort(A[i],n));`

Distribuvana schemata

- Množina uloh je distribuovaná medzi procesmi
- Za behu dochádza k vzájomnému vymenovaniu uloh
- Netrpí nedostatkami spojenými s centralizovaným riešením

Moznosti

- Ako sa určí, kto komu pošle ulohu
- Kto a na základe čoho určí, že je potreba presunúť ulohu
- Koľko uloh má byť presunuté
- Kedy a ako je uloha presunutá

Problem

- Efektivita prenosu ulohy na iný proces

Afinitne planovanie

Vlakna a procesory

- Jednotlivé vlákna sú vykonávané fyzickými procesormi
- Planovanie zaisťuje plánovač OS

Afinitne planovanie(angl. Affinity scheduling)

- Modifikácie algoritmu planovania
- Afinitne planovanie zaisťuje, že výpočetné dávky pridelené jednému procesu/vláknu budú pokiaľ možno pridelené na fyzický tentiež procesor

Vyhody a rizika

- Potencionalne lepsie vyuziti cache
- Striktne lpenie na tomtiez procesore moze narusovat vyvazenost vyuzitia procesorov, teda redukovat vykon aplikacie

Dilema procesu dekompozice-mapovani

Otazka

- Je lepsie najprv dekomponovat na mnoho malych uloh a potom ulohy zhlukovat pri mapovani, alebo naopak obmedzit dekompoziciu aby mapovanie bolo priamociare

Aspekty napomahajuce rozhodnutiu dilematu

- Je cena dekompozice zhodna pre oba scenare?
- Vytvara jemnesia dekompozicia skutocne nezavisle ulohy?
- Je jemnejsia dekompozicia zachovana datova lokalita?
- Je/neni znam pocet jader na cielovej platforme?
- Aka je cena rezie prepinani, najma v situaci kedy pocet vlakien vyrazne prevysuje pocet vypocetnich jadier?

Metody pre redukcie rezie interakcie

Rezia suvisiaca s interakciou

- Rezie suvisiaca s interakciou subeznych uloh je klucovym faktorom ovplyvnujucim vykon paralelnej aplikacie
- Z pohladu rezie interakcie su idealne “embarrassingly parallel” ulohy kde k intrakcii nedochadza

Faktory ovplyvnujuce reziu

- Objemy prenasanych dat
- Frekvencia interakcie
- Cena komunikacie

Zvysovanie datovej locality

Ciel – znizit objem prenasanych dat

- Presun zdielanych datovych struktur do lokalnych kopii
- Minimalizacia objemu zdielanych dat
- Lokalizacia vypoctu(lokalne ukladanie medzivysledkov)
- Rezie protokolov pre udrzanie koherencie lokalnych kopii

Ciel – snazit frekvenciu interakcie

- Priestorova lokalizacia prenasanych dat
- Prenasanie dat a ich okolia(princip cache)
- Viac zprav v jednej(bufferovanie)
- Komunikacia castejsie po mensich kusoch alebo komunikacie menej castejsie po vacsich kusoch

Minimalizacia sucasnych pristupov

Problem – contention

- Pristup k omezenemu zdroju v rovnaky okamih(contention) je rieseny serializaciou poziadavok
- Serializacia poziadavkov sposobuje prodlevy

Mozne riesenie

- Je potreba N subeznych pristupov k datam
- Pristupovane data je mozne rozdelit do N blokov
- A data citat v N po sebe iducich iteraciach
- V kazdej iteracii je kazdy blok citan inym vlaknom
- Cislo citaneho bloku v rte iteracii jtym vlaknom $(r+j)$ modulo N

Prekvryvanie vypoctu s interakciou

Problem

- Cakanie na prijem ci odoslanie dat zposobuje nechcene prodlevy

Vcasne vykonanie akcie – podmienky vykonatelnosti

- Data musia byt vcas pripravené
- Prijmajuca i odosielaajuca strana mozu asynchrone komunikovat
- Existuje dalsia uloha, ktora moze byt riesena po dobu komunikacie

Ine riesenie

- Simulace mechanismu prerusenja(ala operacny system)
- Zadne kvoli rezii sposobene nasilnym riesenim

Replikacia dat ci vypoctu

Problem

- Opakovane drahe pristupy k zdielanym datam

Riesenie pre read-only data

- Pri prvotnej interakcii tvorba kopii dat(datova lokalita)
- Dalej pracovat s lokalnou kopiou
- Zvysuje pametove naroky vypoctu

Riesenie pre read-write data

- Podobne ako v read-only pripade
- Nasobne subezne vypocty toho isteho mozu byt rychlejsie nez citanie a zapis zdielanej hodnoty

Optimalizovane operacie pre kolektivnu komunikaciu

Problem

- Stejna interakcia medzi vsetkymi procesmi vykonava zakladnymi komunikacnymi primitivy je draha

Riesenie – kolektivni komunikacne operacie

- Pre pristup k datam inych vlakien/procesov
- Dolezite pre komunikacne intenzivne vypocty
- Forma efektivnej synchronizacie

Optimalizovane implementacie

- MPI

Prekryvanie interakcii

Problem

- Nedostatocna priepustnost komunikacnej siete, ci absencie kolektivnych komunikacnych operacii

Riesenie

- Zvysit vyuzite komunikacnej siete sucasnou komunikaciou medzi roznyimi parmi procesov

Priklad

- 4 procesy P_1, \dots, P_4
- P_1 chce vsetkym poslat spravu m_1
- $P_1 \rightarrow P_2, P_2 \rightarrow P_3, P_3 \rightarrow P_4$
- $P_1 \rightarrow P_2, (P_2 \rightarrow P_3, P_1 \rightarrow P_4) \rightarrow$ poslu sa dve spravy naraz z P_2 a P_1

Redukcia ceny komunikacie

Komunikacia v nezdielanom adresovom priestore

- Synchronizacia posielanim sprav

- Predavanie dat posielaním zprav

Komunikacia v zdieľanom adresovom priestore

- Synchronizácia korektným prístupom k zdieľaným dátam
- Predavanie dat pomocou zdieľaných dátových štruktúr(FIFO)

Obečné charakteristiky

- **Latence** -> doba potrebná pre doručenie prvého bitu
- **Prenosová rýchlosť** – objem dat prenesených za jednotku času

Latence

- Celková cena pre zahájenie komunikácie – t_s
 - o Čas pre prípravu zpravy/dat
 - o Identifikácia adresata/routovanie
 - o Doba trvania vyliatia informácie z cache do pamäte prípadne na sieťové rozhranie
- Cena "hopov"(preposielaní vo vnútri komunikačnej siete) – t_h
 - o Čas strávený na jednotlivých routeroch v sieti
 - o Doba po ktorú putuje hlavička zpravy z prijímacieho na odosielači port

Prenosová rýchlosť

- Ovlivnené šírkou pásma r
- Cena za prenos jedného slova(Word=2 bajty) -> $t_w = 1/r$

Cena komunikácie

- M -> dĺžka zpravy v slovach
- L – počet liniek cez ktoré zprava putuje
- $T_s + L * (t_h + m t_w)$

Obečné metódy redukcie ceny

- Spojovanie malých správ(amortizuje sa hodnota t_s)
- Kompresia(znížovanie hodnoty m)
- Minimalizácia vzdialenosti(znížovanie hodnoty L)
- Paketovanie(eliminácia rezie spôsobenej jednotlivými hopmi)
 - o $T_s + L * (t_h + m t_w) \rightarrow t_s + L t_h + m t_w$

08- kolektívne komunikačné primitíva

Kvantitatívne parametre komunikácie

Komunikácie(interakcie)

- Predavanie informácií medzi jednotlivými procesmi.

Parametre komunikácie

- **Latencia** -> oneskorenie súvisiace so začiatkom vlastnej komunikácie
- **Sírka pásma(bandwidth)** -> maximálne množstvo dát prenesených za jednotku času
- **Objem** -> množstvo predávaných dát

Cena komunikácie

- T_s -> latencia, t_w -> šírka pásma, m -> objem
- $T_s + m \cdot t_w = \text{cena}$

Príklady

- Za ako dlho napustím hrnček vodou pomocou 2km dlhej záhradnej hadice
 - o Analógia pre vysokú latenciu, aj keď objem dát je malý
- Pri konštantnom čítaní z pamäte trvá získanie kódu inštrukcie a príslušných operandov z pamäte v priemere 5ns, aká je najvyššia reálna rýchlosť vykonávania kódu uloženého v pamäti 4GHz procesorom?
 - o $[1/(5 \cdot 10^{-9}) = 0.2 \cdot 10^9 = 200 \text{ MHz}]$

Efektívnosť interakcie

Pozorovania

- Interakcia jednotlivých úloh/procesov je nevyhnutelná
- Interakcia spôsobuje prodlevy vo výpočte
- Rezie súvisiaca s interakciou by nemala byť dôvodom pre neefektívnosť paralelného spracovania

Záver

- Komunikačné primitívy musia byť čo najefektívnejšie

Topológie komunikačnej siete

Topológie

- Fyzická/logická štruktúra komunikačných kanálov medzi jednotlivými účastníkmi komunikácie

Vlastnosti komunikačnej siete

- Priemer (dĺžka maximálnej najkratšej cesty)
- Konektivita (minimálny počet disjunktných ciest) -> robustnosť
- Stupeň (max počet liniek vychádza z jedného uzla)
- Cena
- Vylúčnosť prístupu (použitie jednou úlohou blokuje ostatné)
- Rozsiritelnosť

- Skalovateľnosť

Sbernice

- Komunikácia cez jedno zdieľané miesto/ zdieľané médium
- Pripustnosť, klesá s počtom uzlov (je treba cache)
- Doležitost: model zdieľaného adresového priestoru

Kličky (úplne siete)

- Private neblokujúce spojenia každého s každým
- Cena: počet liniek je kvadraticky v N
- Cena: stupeň každého uzlu je $N-1$
- Skalovateľné, za podmienky levného zvyšovania stupňa uzlu
- Doležitost: abstraktná predstava siete, logická štruktúra
- Najpoužívanejší model

Prsten (kruh, retaz, 1-rozmerná mriežka)

- Usporiadanie na uzloch
- Privatná neblokujúca komunikácia s 2 najbližšími uzlami
- Doležitost: logická štruktúra komunikácie, pipeline model

Hviezdica

- Spojenie cez jeden centrálny uzol
- Stredový uzol môže byť úzkym miestom
- Lze hierarchicky vrstviť (hviezdica hviezdíc)
- Doležitost: master-slave model

Hyperkostka

- Spojenie n uzlov v tvare $\log_2 n$ -rozmiernej krychle
- Stupeň uzlu je $\log_2 n$
- Priemer siete je $\log_2 n$
- Počet liniek $O(N \cdot \log(N))$
- Postup konštrukcie
 - o Binárne označenie uzlov s identifikátormi 0 až $2^{(\log_2 n)} - 1$
 - o Linka existuje medzi uzly pokiaľ sa označenie líši v jednom bíte
 - o Minimálna vzdialenosť uzlov odpovedá počtu odlišných bitov v označení uzlov

Strom s aktívnymi vnútornými uzly

- Strom, kde účastníkmi komunikácie sú listy i vnútorné uzly
- Kostra hyperkostky – strom s maximálnou hĺbkou $\log_2 n$
- Doležitost: optimálne šírenie informácií

Iné topológie -> veľmi špecifické prípady

- Mriežky

- Cyklické mriežky(torus)

Jeden na vsehny a vsichni na jedneho

One-To-All vysielanie(OTA)

- Uloha posla niekoľko/všetkým ostatným identická data
- Vo výsledku je p kópí originalnej spravy v lokálnych adresových priestoroch adresátu
- Zmena štruktúry dát $m \rightarrow p \cdot m$ (m -je veľkosť spravy)

All-To-One redukcia (ATO)

- Dualná operácia k "One-To-All"
- Niekoľko/všetky ulohy posielajú data jednej ulohy
- Data sa kombinujú pomocou zvoleného asociatívneho operátora
- Vo výsledku je jedna kópia v adresovom priestore cieľovej ulohy
- $M_1 \otimes M_2 \otimes M_3 \dots \otimes M_p \rightarrow M$
- Zmena štruktúry dát $m \cdot p \rightarrow m$
- Dochádza k redukcii

OTA: Pre topologie prsten ci retaz

Naivný spôsob One-To-All pre p procesov

- Poslať $p-1$ správ postupne ostatným procesom
- Úzke miesto: odosielateľ
- Sieť je nevyužitá, komunikuje iba jedna dvojica procesov

Technika rekurzívneho zdvojenia(pripomenutie)

- Najprv prvý proces pošle správu inému procesu
- Potom oba procesy pošlú správu inej dvojici
- Potom štvorice procesov pošlú správu inej štvorici
- Atd
- Prvý proces pošle najviac $\log(p)$ správ
- Súbeh správ na linkách siete
- Optimálna vzdialenosť adresátov pre jednotlivé koly sú $p/2, p/4, p/8, p/16, \dots$

OTA: Pre topologiu d-dimenzionalna mriežka

One-to-All v 2-rozmernej sieti o p uzloch

- Lze chápať ako \sqrt{p} reťazku o \sqrt{p} uzloch
- V prvej fáze sa propaguje informácia do všetkých reťaziek
- V druhej fáze sa zabezpečí propagácia informácie v jednotlivých reťazkách
- Celková cena $2(\log(\sqrt{p}))$

One-To-All v d-rozmernej sieti

- Stejný princíp, veľkosť v jednom rozmere je $p^{1/d}$

- Celkova cena $d(\log(p^{1/d}))$

OTA pre topologiu hyperkostka

Hyperkostka s 2^d vrcholmi

- Aplikuje sa na algoritmus pre sieť vo tvare mriežky
- D-dimenzionalná sieť s hĺbkou siete 2
- Každá fáza prebehne v konštantom čase (poslane 1 spravy)
- Nenastáva subezný prístup na žiadnu z komunikačných liniek
- Celková cena: d

Priklad

- Jak prebehne OneTo-All na hyperkostke s dimenziou 5

ATO: dualne k OTA

All-To-One redukce pre p procesov

- Prebieha dualne k operácii one-to-all

Priklad

- Prvé procesy s lichým ID pošlú správu procesom s ID o jedna menším, kde sa správy zkomponujú s hodnotou, ktorou chcú vyslať procesy so sudým ID
- Nasledne prebehne kombinácia informácií susedných procesov so sudým ID na procesoch, ktorých ID je násobkom 4
- Atd

Univerzálne algoritmy

Pozorovanie

- One to all procedúry sú si podobné
- Jdu nahradiť univerzálnou procedúrou

Univerzálne algoritmy OTA vysielať a ATO redukce

- Predpokladajú 2^d uzlov (hyperkostka)
- Každý uzel identifikovaný bitovým vektorom
- AND a XOR bitové operácie

Cena

- $D(ts+mtw)$
- D počet dĺžok hyperkostky

*

Všetci všem a od všech všem

All-to-All vysielať

- Vsetci posielaju informáciu vsetkym
- Každá uloha posiela jednu zpravu vsetkym ostatnym uloham
- Vo vysledku je p kopii p originalnych zprav v lokalnych adresovych prostoroch adresatu
- Zmena struktury dat $p \times m \rightarrow p \times p \times m$

All-to-All Redukce

- Vsetci posielaju informáciu vsetkym, prichodi spravy sa kombinuju
- Každá uloha ma pre kazdu inu ulohu inu zpravu
- Zmena struktury dat $p \times p \times m \rightarrow p \times m$

Naivne riesenie

- Sekvencne/nasobne pouzitie One-To-All procedur
- Neefektivne vyuzitie siete

ATA: pre topoogiu prsten ci retaz

Prsten

- 1. Faze kazdy uzal posle informáciu svojmu susedovi
- 2 az nta faze: uzly zbieraju a preposielaju prichadzajuce spravy
- Vsetky linky su po celu dobu operacie plne vyuzite
- Optimalny algoritmus

Retez

- Vysielane sprav susedom na obe strany
- Full-duplex linky $\rightarrow n-1$ fazi
- Half-duplex linky $\rightarrow 2(n-1)$ fazi

ATA redukce

- Zpravy po jednej posielane po kruhu tak, ze zprava pre najvzdialenejsi uzal je posлана ako prva
- Pri priechode uzlom sa prikombinuje lokalna zprava pre adresata prave prechadzajucej spravy

ATA: pre topologiu mriezky ci hyperkostky

Mriezka

- 2 faze – \sqrt{p} retiazkov o \sqrt{p} uzloch
- Po prvej fazi uzly maju uzly \sqrt{p} casti z celkoveho poctu p casti zpravy
- Priklad siete 4x4, kazda posiela svoje ID kazdemu

Hyperkostka

- Rozsirenie algoritmu pre siete na d dimenzii
- Po kazdej fazi (z celkoveho poctu d j) je objem zprav zdvojnásobený

ATA redukcia

- Dualny postup
- Po kazdej fazy je objem zprav zredukovany na polovicu

ATA: analyza ceny

Kruh a linearne pole, p uzlov

- $T = (t_s + t_{wm})(p-1)$

2d mriezka, p uzlov

- 1. fáze $(t_s + mt_w)(\sqrt{p}-1)$
- 2. fáze $(t_s + m\sqrt{p}t_w)(\sqrt{p}-1)$
- Celkem: $T = 2t_s(\sqrt{p}-1) + t_{wm}(p-1)$

Hyperkostka, $p=2^d$ uzlov

- $T = \sum_{i=0}^{d-1} (t_s + 2^i t_{wm})$
- $T = t_s \log p + t_{wm}(p-1)$

Pozorovanie

- clen $t_{wm}(p-1)$ se vyskytuje vždy
- pipeline niekoľko OTA operacii

Vsichni vsem individualni komunikacia

All-to-All individualni komunikace (ATA individualni)

- kazda uloha posila rozne data ostatnym uloham
- dojde k vymene 2D pola zprav ($p \times p$) v 1D priestore (p)
- take oznacovane ako totalna vymena
- zmena struktury dat:
 - $p \times (m_1, \dots, m_p) \rightarrow p \times m_1, p \times m_2, \dots, p \times m_p$

Priklad

- transpozice matice ($A^T[i,j] = A[j,i]$)
- matice $n \times n$ mapovana po riadkoch na n uloh
- uloha j ma k dispozicii prvky $[j,0], [j,1], \dots, [j,n-1]$

ATA individualni: pre topologii prsten ci retez

Princip

- najprv ulohy poslu jedným smerom zpravy pre zbyvajucich $p-1$ uloh
- v kazdom dalsom kole kazda uloha vyextrahuje zpravu koja je urcena jej a zbyvajuce spravy preposle
- v poslednom kole vsechny ulohy preposielaju len jednu spravu

Analyza ceny

$$T = \sum_{p-1, i=1} (t_s + t_{wm}(p-i)) = t_s(p-1) + \sum_{p-1, i=1} i t_{wm} = (t_s + t_{wmp}/2)(p-1)$$

ATA individualni pre topologii mriezky

Princip

- Siet je \sqrt{p} retizku o delce \sqrt{p}
- 1. Faze: medzi retizky sa vymenia v jednom smere zpravy tak, aby sa informacia pre kazdy uzol v retazku bola niekde retiazku obsiahnuta
- 2. Faze v ramci retizku sa informacia napropaguje na odpovedajuce miesto

Priklad

- ATA individualni komunikace na sieti 4x4 uzly

Cena

- Každá fáze distribuuje zprávy velikosti $m\sqrt{p}$ mezi \sqrt{p} uzly
- Cena jedné fáze (viz cena pro prsten):
 - o $(t_s + t_{wmp}/2)(\sqrt{p}-1)$
- $T = (2t_s + t_{wmp})(\sqrt{p}-1)$

ATA individualni pre topologii hyperkostky

Princip

- Aplikacia algoritmu pre d-dimenzionalne siete
- Podel jedne dimenzie sa posielajú vždy $p/2$ sprav

Priklad

- Krychle 2x2x2 uzlu

Cena

- V kazdej faze vymenene $mp/2$ dat
- Log p fazi
 - o $T = (t_s + t_{wmp}/2)\log p$
- Navyiac v kazde faze uzly preusporiadavaju/tridia spravy

ATA individualni pre topologii hyperkostky – optimalita

Problem

- Kazdy uzol posielajú a prijímajú $m(p-1)$ dat
- Průměrná vzdálenost, na kterou data putují je $(\log p)/2$

- Celkový provoz na síti je $p \times m(p-1) \times (\log p)/2$
- Počet linek v hyperkostce $\log p/2$
- Optimální algoritmus by měl dosáhnout složitost
 - o $T = t_{wpm}(p-1)(\log p)/2(p \log p)/2 = t_{wm}(p-1)$

Zaver

- Pre ATA individualni komunikaci neni algoritmus pre siete vo tvare mrizi optimalni na sietoch v tvare hyperkostky

*

E-cube routing -> IB109 -> video31 -> este raz pozriet

E-cube routing

- Cesta medzi 2 bodmi je dana poziciou odlisnych bitov tychto bodov, routuje sa od najmenej vyznamneho bytu
- Vyzaduje duplexne linky

Operace kompletnej redukcie a prefixoveho suctu

Operace kompletnej redukcie

Kompletni redukce – All reduce

- Ulohy si vzajomne vymenuju data ktora sa kombinuju
- Lze realizovat jako ATO redukci nasledovanou OTA vysielanim vysledku z predchoczej redukce
- Zmena struktury dat $p \times m \rightarrow p \times m$

Semantika a vyuzite pre ucely synchronizacie

- Uloha nemoze dokoncit operaciu, dokiaľ vsetky participujuci ulohy neprispievaju svojim dielom
- Moze byt pouzite pre realizaciu synchronizacneho primitiva

Implementace

- Naivne pomocou ATO a OTA nebo
- Modifikacia operacie ATA vysielana
 - o Rozdiel od ATA – spravy sa nekumuluju ale kombinuju
- Cena pro hyperkostku je $T = (t_s + t_{wm}) \log p$

Operace prefixoveho suctu

Prefixov sučet

- Ulohy posielaju data ostatnym uloham so stejnym ci mensim ID
- Data sa kombinuju – redukcia
- Zmena struktury dat $p \times m \rightarrow p \times m$

Priklad

- Data v jednotlivych uzloch $\langle 3,1,4,0,2 \rangle$
- Kombinacia pomocou operacie suctu
- Vysledna data $\langle 3,4,8,8,10 \rangle$

Implementacia

- Pouzitie algoritme pre ATA
- Ako sa pozna ze zprava pochadza od uzlu z mensim ID?
- Vsetke posielana data obohatená o ID puvodce

Operace Scatter and Gather

Scatter – tez oznacovane ako “Scan”

- Jedna uloha posielá unikátnu zprávu každej ďalšej ulohy
- One to all individualni (personalized) komunikace
- Na rozdiel od OTA sa žiadne dáta sa neduplikujú
- Zmena štruktúry dát $(m_1, \dots, m_p) \rightarrow m_1, \dots, m_p$

Gather – tez oznacovana ako “Concatenation”

- Jedna uloha zbiera unikátne dáta od ostatných uloh
- All to one individualni (personalized) komunikace
- Na rozdiel od ATO redukcie sa nevyskytuje kombinácia dát
- Zmena štruktúry dát $m_1, \dots, m_p \rightarrow (m_1, \dots, m_p)$

Implementace

- Vuzitie algoritmu pre ATO a OTA
- Celkem $\log p$ fází, s každou fázou sa objem dát zväčšuje
- $T = t_s(\log p) + t_w m(p-1)$

Operace cyklicky posun

Permutace

Permutace

- Obecnější komunikační primitivum
- Súčasne prebiehajúci OTO prerozdelenia dát
- Jedna uloha posielá dáta jednej inej ulohy

Priklad

- Cyklicky posun o q (circular q -shift)
- P uloh
- Uloha i posielá dáta ulohy $(i+q) \bmod p$

Pouzitie

- Špecifické maticové operácie
- Vyhľadávanie vzoru v texte či obraze

Cyklický posun

1-dimenzionalni

- Intuitívne → rotácie o q pozícií
- Smery rotácie závisia na q , určí sa ide výrazu $\min(q, p-q)$

2d siete

- Akcelerácia cyklického posunu s využitím druhej dimenzie než vo ktorej prebieha posun
- Jedna dimenzia má rozmer \sqrt{p} uzlov
- Posun v druhej dimenzii akceleruje o \sqrt{p} krokov

Cena

- Najvzdialenejší posun v jednej dimenzii je $\sqrt{p}/2$
- $T = (t_s + t_w)(\sqrt{p})$

Príklad

- Sieť 4x4 uzlov, cyklický posun o 5

Cyklický posun – sieť v tvare hyperkocky

Princíp

- Mapovanie lineárneho pola na hyperkocku dimenzie d
 - Zrcadený sedý kód (reflected grey code): $I = G(I, d)$
 - $G(0, 1) = 0$
 - $G(1, 1) = 1$
 - $G(i, x+1) = \text{if } i < 2^x G(i, x) \text{ else } 2^x + G(2^{x+1} - 1 - i, x)$
- Q sa vyjadri ako súčet mocnín čísla 2
- Posun prebehne v toľko fázach koľko je členov v súčte

Cena

- Uzly v vzdialenosti mocniny čísla 2, sú od seba vzdialené najviac 2 kroky (vlastnosť kódu)
- Člen v súčte je najviac $\log p$
- Komunikácia prebieha bez blokovania liniek
- $T = (t_s + t_w)(2 \log p)$
- Se zpätným posunom je počet sčítanců $\max. (\log p)/2$
- $T = (t_s + t_w)(\log p)$

09 – programovanie aplikacii pre prostredie s distribuovanou pametou

Principy programovania s predavanim sprav

Vymezeni prostredi

Paradigma – predavanie sprav

- Nezdielany adresovy priestor
- Explicitny paralelizmus

Pozorovanie

- Data explicitne delena a umiestnena do jednotlivych lokalnych adresoych priestorov
- Datova lokalita je klucova vlasnost pre vykon
- Komunikacia vyzaduje aktivnu ucast komunikujucich stran
- Existuju efektivne implementovana podporne knihovne
- Programator je zodpovedny za paralelizaciju algoritmu

Stuktura programov predavanim sprav

Asynchronni paradigma

- Vypocet zacne v rovnaky okamzik(synchrone) ale
- Prebieha asynchrone(Rozne vlakna roznu rychlostou)
- Mozna synchronizacia v jednotlivych bodoch vypoctu
- Neplati "trojuholnikova nerovnost" v komunikacii
 - o Moze sa stat ze ked Vlakno A posielva vlaknam B a C spravy, a vlakno B posielva tu istu spravu C, tak moze sprava od B do C dorazit skor nez sprava od A do C

Dalsie vlastnosti

- Vykazuje nedeterministicke chovanie
- Tazsie prekazanie korektnosti
- Moznost vykonana zcela odlisneho kodu na jednotlivych procesnych jednotkach
- Typicky vsak "single program multiple data"
- Kazda procesna jednotka ma jednoznacu identifikaciju

Send a Receive – zakladne stavebne kamene

```
send(void *sendbuf, int nelems, int dest)
```

`receive(void *recvbuf, int nelems, int src)`

sendbuf

- ukazatel na buffer(blok pameti), kde su umiestnene data pripravene k odoslaniu

recvbuf

- ukazatel na bafr(blok pameti), kam budu umiestnene priate data

nelems

- pocet datovych jednotiek, ktore budu poslane ci priate(delka zpravy)

dest

- adresat odosielanej zpravy, tj ID toho, komu je zprava urcena

src

- odosielateľ spravy, ID toho kto spravu poslal, alebo toho od koho chcem spravu prijať

*

Vystup a efektivita

- asynchrónna semantika `send()` a `receive()`
 - o inak sa `receive()` vola blokujuco
- nutné z dôvodu zachovania výkonu aplikácie

Blokujuce nebafované operácie

Blokujuce operácie

- operácia `send()` ukončená ažtedy, keď je bezpečné vzhľadom k semantike, tj že prijemca odoberie to čo bolo obsahom odosielaného bufferu v okamžiku volania operácie `send()`
- ukončenie operácie `send()` nevyvoláva a negarantuje že prijemca už zprávu prijal
- operácia `receive()` skončí po prijatí dát a jejich umiestnení na správne miesto v pameti

Nebafované operácie

- nie je medzipamät kde by sa dočasne uložil obsah toho čo sa má odoslať
- operácia `send()` skončí až po dokončení operácie komunikácie, tj až keď prijímaci proces zprávu prijme
- v rámci operácií `send()` a `receive()` pred samotným prenosom dát prebieha synchronizácia oboch participujúcich strán(handshake)
- nebafované operácie dávajú najväčší zmysel pri posielaní veľkých dát
- pre malé správy dáva zmysel mať bufované správy

Blokujuce nebafované operácie - problémy

Prodlevy

- sposobene synchronizáciu pred vlastnou komunikáciou
- process ktorý dosiahne bodu, kedy je pripravený komunikovať, čaka až do stejného bodu dospeje i druhý process
- volanie send() a receive() v rovnaký okamžik nelze garantovať na úrovni kodu
- nema vliv, pokiaľ dominuje čas komunikácie

Deadlock

- send a receive môžu vyvolať deadlock ak sa nevolá odpovedajúci receive

*

Blokuje bafrované operácie

Bafr v bafrované komunikácii

- extra pamät zdanlivo mimo adresový priestor procesov
- medzisklad správ pri komunikácii

bafrované komunikacie operácie

- operácia send() skončí v okamihu, kedy odosielané dáta kompletne prekopirované do bafru
- prípadná modifikácia posielaných dát po skončení operácie send(), ale pred začatím vlastnej komunikácie sa neprejaví
- volanie send(), vlastná komunikácia a následne volanie receive() na prijímačej strane sa nemusia časovo prekryvať

Blokuje bafrované operácie- problémy

Režia súvisiaca s bafrovaním

- eliminácia prodlev za cenu rezie bafrovania
- vo vysoko synchronizovaných aplikáciach môže byť horšie než používať blokujúcich nebafrovaných operácií

velkost bafru

- pokiaľ odosielateľ generuje správy rýchlejšie, než je prijímač schopný správy prijímať, veľkosť bafru môže neúmerne rásť (problém producent-konzument)
- pokiaľ je veľkosť pre bafry obmedzená, môže dochádzať (a to nedeterministicky) k situácii, kedy je predom daná veľkosť bafru nedostatečná (buffer overflow)
 - o buď sa správy zahadzujú, alebo sa volanie send() stane nebafrované po nejakej dobe a musí čakať až prijímač správy z bafru zkonsumuje
- prípadne samovoľné blokovanie odosielateľa do tejto doby než odosielateľ prijme nejaké dáta a bafry sa uvoľnia, môže viesť k uviaznutiu, podobne ako v prípade nebafroanej blokujúcej komunikácie

Uviaznutie (i bez blokovania odosielateľa)

- ak najprv zavolam receive a potom send, dostanem deadlock *

asymetricky model blokujuci bafrovane komunikacie

- neexistencia odpovedajucich prostriedkov pre bafrovanu komunikacu na urovni komunikacnej vrstvy
- operacia send() je blokujuca, nebaftrovana
- prijmaci process je prerusený v behu a zprava je prijata do bafu, kde caka dokial prijmaci process nezavola odpovedajucu operaciju receive()
- dedikovane vlakno pre obsluhu komunikacie

neblokujuce operacie

Motivacie

- komunikacia, ktora nespособuje prodlevy

neblokujuce operacie

- volanie funkcie moze skoncit drive nez je to semanticky bezpecne
- existuje funkcie na zistení stavu komunikujucej operacie
- program nesmie modifikovat odosielane data, dokial komunikacia neskonci
- po dobu trvania komunikacie program moze vykonavat kod
- prekryvanie vypoctu a komunikacie

hw realizacia

- DMA

Prehľad komunikacnych modov

	Blokující	Neblokující	
Bafrované	send skončí jakmile jsou data nakopírována do bufru	send skončí jakmile je inicializován DMA přenos	Asynchronní obsluha bufru, vyšší paměťové nároky.
Nebafrované	send skončí až po skončení odpovídajícího receive	send skončí ihned, odešle se pouze požadavek na komunikaci	Nutnost časové synchronizace účastníků komunikace.
	Sémantiku operací nelze porušit	Korektní dokončení operací nutné zjišťovat opakovaným dotazováním se	

Message Passing Interface

Message Passing Interface

Standard MPI

- Standardizuje syntax a sémantiku komunikačních primitiv
- Rozhraní pro C, Fortran
- MPI verze 1.2
- MPI verze 2.0 (paralelní I/O, C++ rozhraní)
- Existují různé implementace standardu
 - o Mpich
 - o LAM/MPI
 - o Open MPI

Minimální použití MPI

MPI funkce	Význam
MPI_Init	Inicializuje MPI
MPI_Finalize	Ukončuje MPI
MPI_Comm_size	Vrací počet participujících procesů
MPI_Comm_rank	Vrací identifikátor volajícího procesu
MPI_Send	Posílá zprávu
MPI_Recv	Přijímá zprávu

Definice typu a konstant: `#include "mpi.h"`

Navratova hodnota pri uspesnom volani funkcie *MPI_SUCCESS*

Kompiplace: *mpicc, mpiCC, mpiC++*

Spustenie programu: *mpirun*

Dobre prostredie pre MPI su zdielane domovske adresare, cez viacero strojov, mat dobre nastavene ssh kluce(nastavenie prihlasovanie sa na stroje bez hesiel)

MPI komunikatory

Komunikacne domeny

- Zdruzovanie participujucich proesorov do skupin
- Skupiny sa mozu prekryat

Komunikatory

- Premenne ktore uchovavaju komunikacne domeny
- Typ *MPI_Comm*
- Su argumentum vsetkych komunikacnych funkcii MPI
- Defaultny komunikator: *MPI_COMM_WORLD*
 - o Komunikator do ktoreho su pripojene vsetky participujuce procesy, taka ko boli pri zapnutí aplikacie

Zistenie velkosti domeny a identifikatoru vramci domeny

- `int MPI_Comm_size(MPI_Comm comm, int *size)`
 - o kolko clenov ma skupina daneho komunikatoru
- `int MPI_Comm_rank(MPI_Comm comm, int *rank)`
 - o zistim rank vramci komunikatoru
- rank je identifikator procesu v danej domene
- rank cislo v intervale `[0,size-1]`

*

Posielanie zprav

```
int MPI_Send(void *buf, int count, MPI_Datatype
datatype, int dest, int tag, MPI_Comm comm)
```

- odosle data odkazovana ukazatelom *buf*
- na data sa pozera ako na sekvenciu instancii typu *datatype*
- odosle sa *count* po sobe iducich instancii
- *dest* je rank adresata v komunikacnej domene urcenej komunikatorom *comm*
- *tag*
 - o prilozena informacia typu *int* v intervale `[0, MPI_TAG_UB]`
 - o pre prijemce viditelna bez citania obsahu zpravy

- o typicky odlisuje typ zpravy

Korespondence datových typu MPI a C

MPI datový typ	Odpovídající datový typ v C
MPI_CHAR	signed char
MPI_SHORT	signed short int
MPI_INT	signed int
MPI_LONG	signed long int
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED_SHORT	unsigned short int
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long int
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double
MPI_BYTE	-
MPI_PACKED	-

MPI ma vlatne typy, lebo umoznuje bezat na roznych typoch procesoroch parallelne -> ine usporiadanie bytov, 32bvs64b

Prijmanie zprav

```
int MPI_Recv(void *buf, int count, MPI_Datatype
datatype, int source, int tag, MPI_Comm comm, MPI_Status
*status)
```

- prijmem zpravu od odosielatela s rankom *source* v komunikacnej domene *comm* s tagom *tag*
- source moze byt *MPI_ANY_SOURCE*
 - o je mi jedno kto to poslal
- tag moze byt *MPI_ANY_TAG*
- Zprava ulozena na adresu urcene ukazatelom *buf*
- Velkost bafru je urcena hodnotami *datatype* a *count*
- Pokial je bafr maly, navratova hodnota bude *MPI_ERR_TRUNCATE*

Prijmanie zprav – MPI_Status

MPI_Status

```
typedef struct MPI_Status {
```

```
int MPI_SOURCE;
int MPI_TAG;
int MPI_ERROR;
};
```

- Vhodne najma v prípade prijmu v rezime *MPI_ANY_SOURCE* alebo *MPI_ANY_TAG*
- *MPI_Status* drži I ďalšie informácie, napríklad skutočný počet prijatých dát (dĺžka správy)

```
int MPI_Get_count(MPI_Status *status, MPI_Datatype
datatype, int *count)
```

Posielanie a Prijímanie správ – semantika

MPI_Recv

- Volanie skončí až sú dáta umiestnené v bufru
- Blokuje receive operácie

MPI_Send

- MPI štandard pripúšťa 2 rôzne zementiky
 - o Volanie skončí až po dokončení odpovedajúcej receive operácie
 - o Volanie skončí, keď sú posielané dáta zkopírované z bufru
- zmena odosielaných dát je vždy semanticky bezpečná
- blokuje send operácie

Možné dôvody uviaznutia

- iné poradie správ pri odosielaní a prijímaní
- cyklické posielanie a prijímanie správ (vecerajúci filozofové → cyklicky si do kruhu posielajú správu)

Súčasné Posielanie a Prijímanie správ

MPI_Sendrecv

- operácia pre súčasné prijímanie a odosielenie správ
- nenastáva cyklické uviaznutie
- buffery pre odosielané a prijímané dáta musia byť rôzne

```
int MPI_Sendrecv (void *sendbuf, int
sendcount, MPI_Datatype senddatatype, int dest, int
sendtag, void *recvbuf, int recvcount, MPI_Datatype
```

```
recvdatatype, int source, int recvtag, MPI_Comm  
comm, MPI_Status *status)
```

Sucasne Posielanie a Prijmanie zprav – Sendrecv_replace

Problemy *MPI_Sendrecv*

- bafry pre odosielana a prijmane data zabieraju 2x tolko miesta
- zlozita manipulacia s datami vďaka 2 rozny**m** bafrom

MPI_Sendrecv_replace

- operace odosle data z bafru a na jejich misto nakopiruje prijata data
- len syntakticky cukor, ono je to stale rozdelen**e** na 2 miesta

```
int MPI_Sendrecv_replace (void *buf, int  
count, MPI_Datatype datatype, int dest, int sendtag, int  
source, int recvtag, MPI_Comm comm, MPI_Status *status)
```

Neblokujuca komunikacia

Neblokující Send a Receive

```
int MPI_Isend(void *buf, int count, MPI_Datatype  
datatype, int dest, int tag, MPI_Comm comm, MPI_Request  
*request)
```

```
int MPI_Irecv(void *buf, int count, MPI_Datatype  
datatype, int source, int tag, MPI_Comm comm, MPI_Request  
*request)
```

- velke I na zaciatku hovor**i** ze sa jedna o neblokujucu funkciu

MPI_Request

- identifikator neblokujucej komunikacnej operacie
- potreba pri dotazovan**i** sa na dokoncenie operacie

Dokoncenie neblokujucich komunikacnych operacii

```
int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status)
```

- Nulový flag znamená, že operace ještě probíhá
 - o 0 prebieha, nenula neprebieha/bola dokon**c**ena

- Při prvním volání po dokončení operace se naplní status, zničí request a flag nastaví nature

`int MPI_Wait(MPI_Request *request, MPI_Status *status)`

- Blokující čekání na dokončení operace
 - o Urobi z neblokujucej operacie blokujucu
- Po dokončení je request zničen a status naplněn

`int MPI_Request_free(MPI_Request *request)`

- Explicitní zničení objektu request
- Nemá vliv na probíhající operaci

Neblokujuce operacie – poznámky

Parovanie

- Neblokujuce a blokujuci send a receive sa mozu lubovolne parovat

Uviaznutie

- Neblokujuce operacie riesia vacsinu problem s uviaznutim
- Neblokujuce operacie maju vzdy vyssie pametove naroky
- Pozdejsie zahajena neblokujuca operacia moze skoncit drive

Kolektivni komunikace

Kolektivni operace

- Mnozina participujucich procesov je urcena komunikacnych domenou(`MPI_Comm`)
- ***Vsetky procesy v domene musia volat odpovedajucu MPI funkciu***
- Obecne sa kolektivne operacie nechovaju ako bariery, tj jeden process moze dokoncit volani funkcie drive, nez jiny process vobec dosiahne miesta volania kolektivnej operacie
- Forma virtualnej synchronizacie
- Nepouzivaju *tagy* (vsetci vedia aka operacia sa vykonava)
- Pokial je nutne specifikovat zdrojovy, ci cielovy process, musia tak urobit vsetky participujuce procesy a jejich volba cieloveho ci zdrojoveho procesu musia byt zhodna
- MPI podporuje 2 varianty kolektivnych operacii
 - o Posielaju sa rovnake velke data (napr `MPI_Scatter`) -> nevektorove
 - o Posielaju sa rozne velke data (napríklad `MPI_Scatterv`) -> vektorove
- `MPI_Scatter()` -> distribuuje mnozinu sprav do zdielaneho prostredia/siete

Bariera

`int MPI_Barrier(MPI_Comm comm)`

- Realizuje bareru

- Volanie funkcie skonci pokiaľ všetky zúčastnené procesy zavolajú MPI_Barrier
- Není nutné aby volaná funkcia bola na "rovnakom mieste v programe"

Kolektivní komunikací primitive a MPI

Operace	Jméno MPI funkce
OTA vysílání	MPI_Bcast
ATO redukce	MPI_Reduce
ATA vysílání	MPI_Allgather
ATA redukce	MPI_Reduce_scatter
Kompletní redukce	MPI_Allreduce
Gather	MPI_Gather
Scatter	MPI_Scatter
ATA zosobněná komunikace	MPI_Alltoall

Vsesmerové vysielanie – OTA

```
int MPI_Bcast(void *buf, int count, MPI_Datatype
datatype, int source, MPI_Comm comm)
```

- Rozosielať dáta uložené v bafre *buf* procesu *source* ostatným procesom v domene *comm*
- Kromě *buf* musia byť parameter funkcie zhodné vo všetkých participujúcich procesoch
- Parameter *buf* na ostatných procesoch slúži pre identifikáciu bafre pre príjem dát

Redukce

```
int MPI_Reduce(void *sendbuf, void *recvbuf, int count,
MPI_Datatype datatype, MPI_Op op, int target, MPI_Comm
comm)
```

- Dáta ze *sendbuf* zkombinované operáciou *op* do *recvbuf* procesu *target*
- Všetci zúčastnení musia poskytnúť *recvbuf* I keď výsledok uložený iba na procese *target*
- Hodnoty *count*, *datatype*, *op*, *target* musia byť zhodné vo všetkých volajúcich procesoch
- Možnosť definovať vlastnú operáciu typu *MPI_Op*

Operatory redukce

Operace	Význam	Datové typy
MPI_MAX	Maximum	C integers and floating points
MPI_MIN	Minimum	C integers and floating points
MPI_SUM	Součet	C integers and floating points
MPI_PROD	Součin	C integers and floating points
MPI LAND	Logické AND	C integers
MPI_BAND	Bitové AND	C integers and byte
MPI_LOR	Logické OR	C integers
MPI BOR	Bitové OR	C integers and byte
MPI_LXOR	Logické XOR	C integers
MPI_BXOR	Bitové XOR	C integers and byte
MPI_MAXLOC	Maximum a minimální pozice s maximem	Datové dvojice
MPI_MINLOC	Minimum a minimální pozice s minimem	Datové dvojice

MAXLOC -> prvá složka dvojice je maxlock samotný, druhé je ID procesu na kterém to maximum bylo

MINLOC analogicky

Datové páry a kompletne redukcie

MPI datové páry	C datový typ
MPI_2INT	int, int
MPI_SHORT_INT	short, int
MPI_LONG_INT	long, int
MPI_LONG_DOUBLE_INT	long double, int
MPI_FLOAT_INT	float, int
MPI_DOUBLE_INT	double, int

- Druhým argumentum je int pretože *rank* je vždy int

```
int MPI_Allreduce(void *sendbuf, void *recvbuf, int  
count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

Prefixovy sucet

```
int MPI_Scan(void *sendbuf, void *recvbuf, int count,  
MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

- Vykonáva prefixovú redukciu

- Process s rankom **I** ma vo výsledku hodnotu vzniknu redukcii procesov s rankom **0 az I** včetne
- Inak zhodne s redukcii

```
int MPI_Gather(void *sendbuf, int
sendcount,MPI_Datatype senddatatype,void *recvbuf, int
recvcount,MPI_Datatype recvdatatype,int target,
MPI_Comm comm)
```

- Vsetci posielaju stejný typ dat
- Cieľový process obdrží **p** bajtov zoradených podľa ranku odosielateľa
- *Recvbuf, recvcount, recvbuf* platné len pre process s rankom target
- *Recvcount* je počet odoslaných dat jedným procesom, nikoliv celkový počet prijímacích dat

Allgather

```
int MPI_Allgather(void *sendbuf, int
sendcount,MPI_Datatype senddatatype,void *recvbuf, int
recvcount,MPI_Datatype recvdatatype,MPI_Comm comm)
```

- Bez určenia cieľového procesu, výsledok obdržia všetci
- *Recvbuf, recvcount, recvdatatype* musia byť platné pre všetky volajúce procesy

Vektorová varianta Gather

```
int MPI_Gatherv(void *sendbuf, int
sendcount,MPI_Datatype senddatatype,void *recvbuf, int
*recvcounts,int *displs,MPI_Datatype recvdatatype,int
target, MPI_Comm comm)
```

- Odosielatelia môžu odosielať rôzne veľké dáta (Rôzne hodnoty sendcount)
- Pole *recvcounts* udáva, koľko dat môže byť prijaté od jednotlivých procesov
- Pole *displs* udáva kde v bajtovi *recvbuf* začínajú dáta jednotlivých procesov
 - o 0412 -> 0 rank má 4 bajty, process s rankom 1 8 bajtov, 12 hovorí kam skocit

Scatter

```
int MPI_Scatter(void *sendbuf, int
sendcount,MPI_Datatype senddatatype,void *recvbuf, int
recvcount,MPI_Datatype recvdatatype,int source,
MPI_Comm comm)
```

- Posielajú rôzne dáta STEJNEJ veľkosti všetkým procesom


```
int MPI_Scatterv(void *sendbuf, int *sendcounts, int
*displs, MPI_Datatype senddatatype, void *recvbuf, int
recvcounts, MPI_Datatype recvdatatype, int source,
MPI_Comm comm)
```

- Posielaju rozne data ROZNEJ velkosti vsetkym procesom

Alltoall

```
int MPI_Alltoall(void *sendbuf, int
sendcount, MPI_Datatype senddatatype, void *recvbuf, int
recvcount, MPI_Datatype recvdatatype, MPI_Comm comm)
```

- Posiela stejne velke casti bafu sendbuf jednotlivym procesom v domene comm
- Process **i** obdrzi cast velkost sendcount ktora zacina na pozicii sendcount*i
- Kazdy process ma v bafu recvbuf na pozici recvcount * I data velkosti recvcount od procesu i

Vektorova varianta Alltoall

```
int MPI_Alltoallv(void *sendbuf, int *sendcounts, int
*sdispls, MPI_Datatype senddatatype, void *recvbuf, int
*recvcounts, int *rdispls, MPI_Datatype
recvdatatype, MPI_Comm comm)
```

- Pole sdispl urcuje, kde začínajú v bafu sendbuf data určená jednotlivým procesom.
- Pole sendcounts urcuje množství dat odesílaných jednotlivým procesom.
- Pole rdispls a recvcounts udávajú stejné informace pro přijatá data.

Skupiny, komunikatory a topologie

```
int MPI_Comm_split(MPI_Comm comm, int color, int
key, MPI_Comm *newcomm)
```

- Kolektivni operace, musi byt volana vsemi
- Parameter *colour* urcuje vyslednu skupinu/domenu
- Parameter *key* urcuje rank vo vyslednej skupine
 - o Pri zhode *key* rozhoduje povodny rank

- Stvrty argument je navratova hodnota pre novy argument pre danu domenu(biele/cervine/etc)

Mapovani procesu

Nevyhody "rucneho" mapovania

- Pravidla mapovania urcene v dobe kompilacie programu
- Nemusia odpovedat optimalnemu mapovaniu
- Nevhodne najma v pripadoch nehomogenneho prostredia

Mapovanie cez MPI

- Mapovanie urcene za behu programu
- MPI knihovna ma k dispozicii(aspon ciastocnu) informaciu o sietovom prostredii(napriklad pocet pouzitych procesorov v jednotlivych participujucich uzloch)
- Mapovanie navrhnuté s ohľadom na minimalizáciu ceny komunikácie

Kartezske topologie

```
int MPI_Cart_create (MPI_Comm comm_old, int ndims, int
*dims,int *periods, int reorder, MPI_Comm *comm_cart)
```

- V tejto funkcii sa specifikuje na kolko skupin/ako velkych skupin sa maju procesy rozdelit a MPI sa posnazi najst take mapovanie aby vyhovelo a minimalizovalo ceny komunikacii atd
- Pokial je v povodnej domene *comm_old* dostacocny pocet procesorov, tak vytvoru novu domenu *comm_carts* virtualni kartezskou topologiou
- Funkciu musia zavolat vsetky procesy z domeny *comm_old*
- Parametry kartezskej topologie
 - o Ndsims – pocet dimenzii
 - o Dims[] – pole rozmerov jednotlivych dimenzii
 - o Periods[] – pole priznakov cyklickej uzavrenosti dimenzii
- priznak *reorder* znaci ze ranky procesov sa maju v ramci novej domeny vhodne preusporiadat
- nepouzite procesy oznacene rankem *MPI_COMM_NULL*

Koordinaty procesorov v Kartezske topologiech

- komunikacne primitive vyzaduju rank adresata
- preklad z koordinatu(cords[]) do ranku

```
int MPI_Cart_rank (MPI_Comm comm_cart,int *coords, int
*rank)
```

- preklad z rank una koordinaty
- maxdims je velkost vstupneho pola coords[]

```
int MPI_Cart_coord(MPI_Comm comm_cart, int rank, int
maxdims, int *coords)
```

Delenie kartezyckych topologii

```
int MPI_Cart_sub(MPI_Comm comm_cart, int
*keep_dims, MPI_Comm *comm_subcart)
```

- pre delenie kartezyckych topologii na topologie s mensou dimenziou
- pole priznakov *keep_dim* urcuje, ci bude odpovedajuca dimenze zachovana v novom deleni

Priklad

- topologie o rozmeroch 2x4x7
- hodnotou *keep_dims*={true,false,true}
- vzniknu 4 nove domeny o rozmeroch 2x7

Pripadova studie implementace verifikacniho nastroje DiVinE

Divine-cluster

- softwarovy anstroj pre verifikaciu protokolov(LTL MC)
- problem detekcie akceptujuciho cyklu v grafu
- algoritmy paralelne prochazi graf konecneho automatu
- standardni pruzkumova dekompozice

Algoritmus MAP

- detekuje zda existuje akceptujuci vrchol, který je svůj vlastní predchudze

Algoritmus OWCTY

- oznacuje vrcholy, které nejsou sucastou akceptujuceho cyklu
 - o nemaju priamych predchodcov – nelezi na cyklu
 - o nema akceptujucich predkv – nelezi na akceptujucom cycle

10 – Analytický model paralelných programů

Vyhodnocování sekvencných algoritmov

- casova a priestorova zlozitost
- funkcie a velkost vstupu

vyhodnocovanie paralelnych algoritmov

- paralelny system = algoritmus + platforma
- slozitost
- prinos paralelizmu
- prenositelnost
- atd

Rezie(overhead) paralelnych programov

N-nasobne zrychlenie

- s pouzitim n-nasobnych HW zdrojov, lze ocekavat n-nasobne zychlenie vypoctu
- prakticky nastava zriedka z dovodu rezii paralelného riesenie
- pokial nastane, je jeho dovodom casto neoptimalne rieseni sekvenčného algoritmu(moze byt aj viac nez n-nasobne zrychlenie pre n-nasobne zdroje)

Duvody

- rezie – interakce, prostoje, zlozitost paralelného algoritmu
- nerovnomerne zvysovanie HW zdrojov
 - o zvyšenie počtu procesorov nepomože, pokial aplikacia neni zavisla na čistom výpočetnom výkone

Komunikace

- cena samotnej komunikacie
- priprava dat k odoslaniu, zpracovanie prijatych dat

Prostojel(lkovani, Idling)

- nerovnomerna zataz na jednotlivé výpočetné jadra
- cakanie na zdroje ci data
- nutnost synchronizacie asynchronných výpočtov

vacsia výpočetna zlozitost paralelného algoritmu

- sekvenčný algoritmus nejde paralelizovat(DFS postorder)
- typicky existuje viacero paralelných algoritmov, je nutne charakterizovat, cim sa plati za paralelizmus

Metriky vykonnosti

Cas vypoctu – execution time

Otazky

- ako merit vykon/kvalitu paralelného algoritmu
- podľa čoho určiť najlepší/najvhodnejší algoritmus pre danú platformu

Cas vypoctu – sekvencny algoritmus

- doba, ktora uplynie od spustenia vypoctu do jeho ukoncenia
- T_S

Cas vypoctu – paralelny algoritmus

- Doba, ktora uplynie od spustenia vypoctu do doby, kedy skonci posledny z paralelnych vypoctov
 - o Zapocitava sa aj cas napríklad distribuovania dat medzi participujucich procesov ci zbierania vysledkov, ak sa vysledok uklada do suboru niekde na server, rata sa aj to
- Zahnuje distribuciu vstupnych i zber vystupnych dat
- T_P

Celkova rezia paralelizmu (Total overhead)

Celkova rezia

- Oznacujeme TO
- Veskere element sposobujuce reziu paralelného riesenia
- Celkovy cas paralelného vypoctu bez casu potrebného pre vypocet problem optimalnym sekvencnym riesenim
- Paralelne a sekvencne algoritmy mozu byt zcela odlisne
- Sekvencnych algoritmu moze existovat viac
- T_S cas vypoctu najlepsiho sekvencneho algoritmu (riesenie)

Funkcia celkovej rezie

- Jaky je cas vypoctu jednotlivych procesov
- Doba po skonceni vypoctu jedneho procesu do skonceni celeho paralelného vypoctu sa povazuje za reziu (idling)
- $TO = pT_P - T_S$
 - o p je pocet vlakien, T_P paralelny cas, T_S optimalny sekvencny cas

zrychlenie (speed-up)

Zakladny prinos paralelizacie

- problem idu vzdy riesit sekvencnym algoritmom
- paralelizaciou lze dosiahnut iba zrychlenie vypoctu
- ostatne vyhody su diskutabilne a tazko meritelne

Zakladna miera ucinnosti samotnej paralelizacie

- pomer casu potrebnych k vyrieseniu ulohy na jednej procesnej jednotke a p procesnych jednotkach
- uvazuje sa vzdy cas najlepsiho sekvencneho algoritmu

- v praxi sa často (a nesprávne) používa čas potrebný k vyriešeniu úlohy paralelným algoritmom spusteným na jednej procesnej jednotke
- $S = TS / TP$

Teoretická hranica zrýchlenia

- S použitím p procesných jednotiek je maximálne zrýchlenie p

Super-lineárne zrýchlenie

- Jav, keď zrýchlenie je väčšie ako p

Pozorovanie

- Zrýchlenie lze meriť asymptoticky

Příklad – sčítaní n čísel s použitím n procesných jednotiek

- Sekvenčne je potreba $\Theta(n)$ operácií, čas $\Theta(n)$
- Paralelne je potreba $\Theta(n)$ operácií, ale v čase $\Theta(\log n)$
- Zrýchlenie $S = \Theta(n / \log n)$

Super-lineárne zrýchlenie

Falošné super-lineárne zrýchlenie

- Uvažme datovú distribúciu na 2 procesné jednotky a operáciu kvadratickú vo veľkosti dát
- Zrýchlenie $S = n^2 / ((n/2)^2)$ pri použití 2 procesných jednotiek
- Problém – neuvažovaný optimálny sekvencný algoritmus

Skutočné super-lineárne zrýchlenie

- Väčšia agregovaná veľkosť cache pamäte
- Pri znížení množstva dát na jednu procesnú uzol sa účinnosť cache pamäte zvýši
- Výpočet je na $1/p$ dátach viac ako p -krát rýchlejší, pretože proces nemusí stále siahť do pamäte

Super-lineárne zrýchlenie závislosti na instancii problému

- Priezkumová dekompozícia úlohy pri hľadaní riešenia
 - Superlineárne zrýchlenie je priezkumovej dekompozícii vlastné
- Paralelná verzia môže vykonať menšie množstvo práce
- V konkrétnej instancii lze paralelne prehľadávanie simulovať i sekvencným algoritmom, občas ale nelze

Amdahlov zákon

Otázka

- Aké je najväčšie možné zrýchlenie systému, pokiaľ sa paralelizáciou urychlí iba čas výpočtu

Amdahluv zákon

- $S_{max} = 1 / ((1-P) + (P/S_p))$
- P – podíl systému, který je urychlen paralelizací
- S_p – zrychlení dosažené paralelizací nad daným podílem

Příklad

- Paralelizací lze urychlit 4-násobně 30% kódu, tj. $P = 0.3$ a $S_p = 4$.
- Maximální celkové zrychlení S_{max} je

$$S_{max} = \frac{1}{(1 - 0.3) + \frac{0.3}{4}} = \frac{1}{0.7 + 0.075} = \frac{1}{0.775} = 1.29$$

Efektivita

Fakta

- Pouze ideální paralelní systém s p procesními jednotkami může dosahovat zrychlení p
- Část výpočtu vykonávána na jednom procesoru je spotřebována rezíou. Procesor nevenuje 100% výkonu řešení problému

Efektivita

- Lze definovat jako podíl zrychlení S vůči počtu jednotek p
- $E = S/p = (TS/TP)/p = TS/(p \cdot T_p)$
- Zrychlení je v praxi $< p$, efektivita v rozmezí $(0,1)$
- “podíl času, po který jednotka vykonává užitečný kód”
- V odborné komunitě nejvíce používáno pro porovnávání kvality paralelních algoritmů

Příklad 1

- Jaká je efektivita sčítání n čísel na n procesorech?
- $S = \Theta\left(\frac{n}{\log n}\right)$
- $E = \Theta\left(\frac{S}{n}\right) = \Theta\left(\frac{1}{\log n}\right)$

Příklad 2

- Na kolik se zkrátí 100 vteřinový výpočet při použití 12 procesorů při 60% efektivitě?
- $E = \frac{T_s}{pT_p}$
- $0.6 = \frac{100}{12x} \implies x = \frac{100}{0.6 \cdot 12} \implies x = 13.88$

Cena

Cena řešení problem na danom paralelnom system

- Sucin poctu procesnych jednotiek a doby paralelného výpočtu:
 - $C = p \times T_p$
- Označovaná také “množstvo práce”, ktoré paralelný system vykona, nebo taktiez “ pT_p product”
- Alternativne lze použít pre výpočet účinnosti ($E = T_s/C$)

Cenovo optimalny paralelny system

- Cena sekvencného výpočtu odpoveda najlepšiemu T_s
- Paralelný system je cenovo optimalny pokiaľ cena riešenia rastie asymptoticky stejne rychle ako cena sekvencného výpočtu
- Cenovo optimalny system musi mat efektivity rovnú $O(1)$

Příklad

- Scítanie n čísel na n procesoroch
- $C = O(n \cdot \log n)$, není cenovo optimalny
- Uvažme n procesorový systém, který třídí (řadí) seznam n čísel v čase $(\log n)^2$
- Nejlepší sekvencní algoritmus má $T_s = n \log n$
- $S = \frac{n}{\log n}$, $E = \frac{1}{\log n}$
- $C = n(\log n)^2$
- Není cenově optimální, ale pouze o faktor $\log n$
- Uvažme, že v praxi $p \ll n$ (p je mnohem menší než n)
- $T_p = n(\log n)^2 / p$ a tedy $S = \frac{p}{\log n}$
- Konkrétně: $n = 2^{10}$, $\log n = 10$ a $p = 32 \Rightarrow S = 3.2$
- Konkrétně: $n = 2^{20}$, $\log n = 20$ a $p = 32 \Rightarrow S = 1.6$
- Závěr: cenová optimalita je nutná

Vliv granularity na cenovu optimalitu

Tvrzení

- Volbou granularity lze ovlivnit cenu paralelného řešení

Pozorování

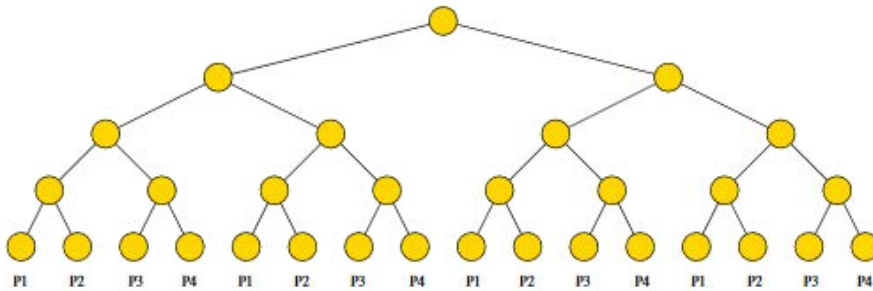
- V praxi $p \ll n$ přesto navrhujeme algoritmy tak, aby granularita byla až na úrovni jednotlivých položek vstupu, tj. předpokládáme že $p = n$

Deskalování(scale-down)

- Umele snižujeme granularitu(vytváříme hrubší úlohy)
- Snižujeme overhead spojený s komunikací
- Může ovlivnit cenu a cenovou optimalitu

Úkol

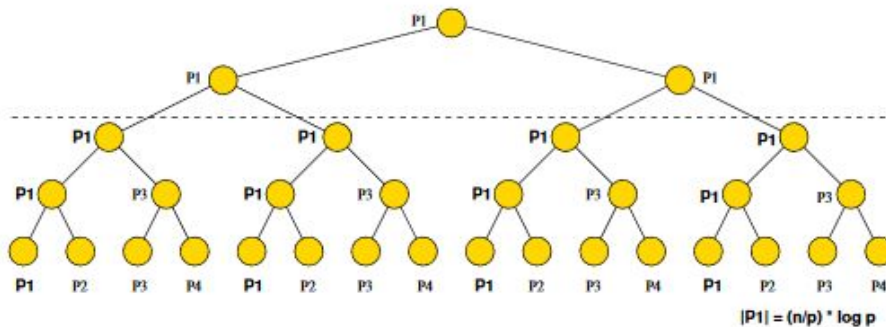
- Sčítání n čísel na p procesorech ($n = 16$, $p = 4$)
- Algoritmus „stromového“ sčítání.
- Mapování ($i \bmod p$), tj. (n/p) čísel na 1 procesor.



Vplyv granularity na cenovu optimalitu

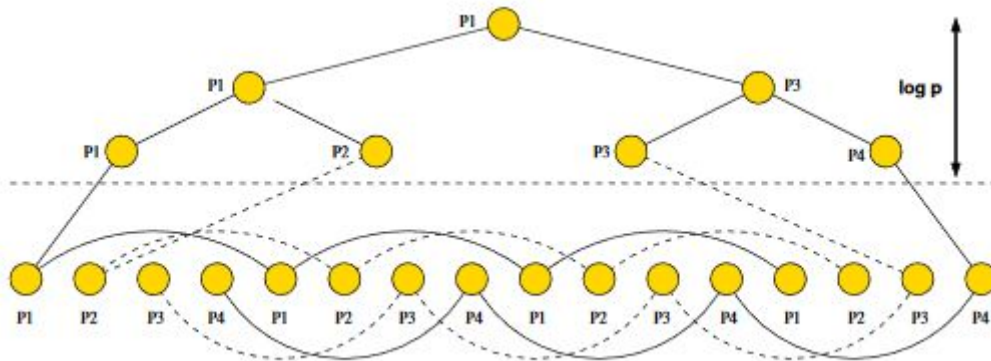
Analýza ceny

- První část výpočtu proběhne v čase $\Theta((n/p)\log p)$.
- Zbylá část výpočtu probíhá lokálně v paměti jednoho procesoru v čase $\Theta(n/p)$.
- Celkový T_P je $\Theta((n/p)\log p + n/p) = \Theta((n/p)\log p)$
- Cena $C = p * T_P = \Theta(n\log p)$, tj. není cenově optimální



Deškálovaný algoritmus a analýza jeho ceny

- Každá jednotka nejprve sečte data lokálně v čase $\Theta(n/p)$
- Problém redukován na sčítání p čísel na p procesorech
- Celkový T_P je $\Theta(n/p + \log p)$
- Cena $C = \Theta(n + p\log p)$,
- Cenově optimální, pokud $n > p\log p$ (cena $C = \Theta(n)$)



Skalovatelnost paralelních programů

Skalovatelnost

- Zachování výkonu a efektivity při zvyšování počtu procesních jednotek a zvacsování objemu vstupních dat

Pozorovanie

- Programy su testovane na malých vstupných detaoch na systemoch s malým počtom procesných jednotiek

- Algoritmus ktorý vykazuje najlepší výkon na testovaných dátach, sa môže ukázať byť tým najhorším algoritmom pri použití na skutočných dátach
- Odhad výkonu aplikácie nad realnými vstupnými dátami a väčším počtom procesných jednotiek je komplikovaný

Záver

- Je treba uvažovať analytické techniky pre vyhodnotovanie výkonu a škálovania

Charakteristiky škálovateľnosti - efektívnosť

Efektívnosť paralelných programov:

$$E = \frac{S}{p} = \frac{T_S}{pT_P} = \frac{1}{1 + \frac{T_O}{T_S}}$$

Celková rezia (TO)

- Prítomnosť sekvencnej časti kódu je nevyhnutelná. Pre jej vykonanie je potrebná časť serial. Po túto dobu sú ostatné procesné jednotky nevyužívané
- $TO > (p-1)t_{serial}$
- TO rastie minimálne lineárne vzhľadom k p , často aj asymptoticky viac

Úloha konštantnej veľkosti (TS fixná)

- So zvyšujúcim sa p , efektívnosť nevyhnutne klesá
- **Nevyhnutný podiel všetkých paralelných programov**

Problém

- Úkol: Sečíst n čísel s využitím p procesních jednotek.
- Uvažme cenově optimální verzi algoritmu.
- $T_p = (\frac{n}{p} + \log p)$.
- Lokální operace stojí 1, komunikace 2 jednotky času.

Charakteristiky řešení

- $T_p = \frac{n}{p} + 2\log p$
- $S = \frac{n}{\frac{n}{p} + 2\log p}$
- $E = \frac{1}{1 + \frac{2p\log p}{n}}$

Charakteristiky škálovateľnosti – škálovateľný systém

Pri zachovaní T_s

- Rastuci pocet procesnych jednotiek znizuje efektivitu

Pri zachovaní počtu procesných jednotiek

- Rastuci T_S (veľkosť vstupných dát) má tendenciu zvyšovať efektivitu

Skalovateľné systémy

- Efektivitu lze zachovať pri subeznom zvyšovaní T_S a p
- Zaujímame sa o pomer, v ktorom sa T_S a p musia zvyšovať, aby sa zachovala efektivita
- Čím menší je pomer T_S/p , tým lepšia skalovateľnosť

Izoefektivita ako metrika skalovania

Vztahy

- $T_P = \frac{T_S + T_O(W, p)}{p}$ W – veľkosť vstupných dát
- $S = \frac{T_S}{T_P} = \frac{pT_S}{T_S + T_O(W, p)}$
- $E = \frac{S}{p} = \frac{T_S}{T_S + T_O(W, p)} = \frac{1}{1 + T_O(W, p)/T_S}$

Konstantní efektivita

- Efektivita je konstantní, pouze pokud $T_O(W, p)/T_S$ je konstantní.
- Úpravou vztahu pro efektivitu: $T_S = \frac{E}{1-E} T_O(W, p)$.
- Při zachování míry efektivity lze $E/(1-E)$ označit jako konstantu K .

Funkce izoefektivity (stejná efektivita)

$$T_S = K T_O(W, p)$$

Pozorování

- Pri konstantnej efektivitě lze T_S vyjádřit jako funkci p
- Vyjadruje vzťah, ako sa musí zvýšiť T_S pri zvýšení p
- Nízka funkcia hovorí, že systém je snáze skalovateľný
- Izoefektivita nelze měřit u systémů, které neskálují

Příklad

- Součet n čísel na p procesorech
- $E = \frac{1}{1+(2p \log p)/n} = \frac{1}{1+T_O(W,p)/T_S}$
- $T_O = 2p \log p$
- Funkce izoefektivity: $T_S = K2p \log p$
- Funkce izoefektivity: $\Theta(p \log p)$

- Při zvýšení procesních jednotek z p na p' se pro zachování efektivity musí zvětšit velikost problému o faktor $(p' \log p')/(p \log p)$.

Cenova optimalita a izoefektivita

Optimalni cen

- Pre cenove optimalne systemy pozadujeme $pTP = O(TS)$
- Po dosadeni zo zakladnych izo vztahov dostavame
 - o $TS + TO(W,p) = O(TS)$

Cenova optimalita moze byt zachovana pouze ak:

- Rezie je najvyse linearni voci zlozitosti sekvenecneho algoritmu, tj funkce $TO(W,p) \in O(Ts)$
- Slozitosť sekvenecneho algoritmu je minimalne tak velka ako rezie, tj funkce $Ts \in \Omega(To(mW,p))$

Spodni odhad na izoefektivitu

Izoefektivita

- Cim nizsi/mensi funkce, tim lepsi skalovatelnost
- Snaha o minimalnu hodnotu izofunkce

Sublinearni izoefektivita

- Pre problem tvoren n jednotkami prace, je otpimalna cena dosiahnute pre maximalne N procesnych jednotiek
- Pridavam dalsich jednotiek vyusti v idling niektorych jednotiek a znizovanie efektivity
- Aby sa efektivita nesnizovala, tak musi mnozstvo prace rast minimalne linearne vzhľadom k p
- Funkce izoefektivity nemuse byt sublinearni

Minimalni doba vypoctu

Otázka

- Jaká je minimální možná doba výpočtu (T_P^{min}) při dostupnost dostatečného počtu procesních jednotek?

Pozorování

- Při rostoucím p se T_P asymptoticky blíží k T_P^{min} .
- Po dosažení T_P^{min} se T_P zvětšuje spolu s p .

Jak zjistit T_P^{min} ?

- Nechť p_0 je kořenem diferenciální rovnice $\frac{dT_P}{dp}$.
- Dosazení p_0 do vztahu pro T_P dává hodnotu T_P^{min} .

Příklad – součet n čísel p procesními jednotkami

- $T_P = \frac{n}{p} + 2\log p$, $\frac{dT_P}{dp} = -\frac{n}{p^2} + \frac{2}{p}$
- $p_0 = \frac{n}{2}$, $T_P^{min} = 2\log n = \Theta(\log n)$

Asymptotická analýza paralelních programů

Algoritmus	A1	A2	A3	A4
p	n^2	$\log n$	n	\sqrt{n}
T_P	1	n	\sqrt{n}	$\sqrt{n} \log n$
S	$n \log n$	$\log n$	$\sqrt{n} \log n$	\sqrt{n}
E	$\frac{\log n}{n}$	1	$\frac{\log n}{\sqrt{n}}$	1
pT_P	n^2	$n \log n$	$n^{1.5}$	$n \log n$

Otázky

- Jaký algoritmus je nejlepší?
- Použily jsme vhodné odhady složitosti?
- Je dostupná odpovídající paralelní architektura?