

Principer för design av parallella algoritmer

Grama et al.
Introduction to Parallel Computing
Kapitel 3
Robert Granat

Grundläggande om parallell algoritmdesign

- 1 **Seriell algoritmdesign** - sekvens av steg för att lösa givet problem m h a e n dator, ett slags "recept"
- 1 Icketrivial **parallell algoritmdesign** - utökar seriella designen:
 - **Identifiera** delar av arbetet som kan utföras parallellt
 - **Mappa** parallella delar på flera processer som arbetar parallellt
 - **Distribuera** indata, resultat och mellanliggande värden
 - **Hantera** återkomst till gemensamt eller delat data
 - **Synkronisera** processer under programexekveringens gång
- 1 Framför allt:
 - **Dela upp arbetet** i mindre delar
 - **Tildela dessa delar** till olika processorer
 - Annat kan falla sig naturligt längs vägen...

Begrepp för denna föreläsning

- 1 Dekomposition
- 1 Uppgifter (bättre: eng. *tasks*)
- 1 Beroendegrafer
- 1 Mappning
- 1 Metoder för att dölja interaktion
- 1 Modeller för parallella algoritmer

Detaljerad översikt

Dekomposition (uppdelning)

- 1 Rekursion
- 1 Data
- 1 Undersökande
- 1 Spekulativ
- 1 Hybrid

Arbetsuppgifter (tasks)

- 1 Egenskaper
- 1 Egenskaper för interaktion mellan uppgifter

Beroendegrafer

Mappning

- 1 Statisk
- 1 Dynamisk

Metoder för att dölja interaktion

- 1 Maximera lokalitet
- 1 Minimera flaskhalsar
- 1 Överlappa beräkningar och kommunikation
- 1 Kollektiva kommunikationsoperationer
- 1 Överlappa interaktioner
- 1 Replikera data
- 1 Extra beräkningar

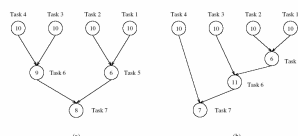
Modeller för parallella algoritmer

- 1 Dataparallell
- 1 Uppgiftsgraf
- 1 Arbetspool
- 1 Master-slave
- 1 Pipe-line

Dekomposition – olika typer av uppdelning

- 1 Lösa ett problem **parallellt** \Rightarrow **dela upp** beräkningarna
- 1 Uppdelningen definierar **uppgifterna**
- 1 Uppdelningen definierar DAG-beroendegrafen (riktad, acyklisk, noder=tasks, kanter=beroenden)
- 1 Första (och viktigaste) steget i designen av en parallell algoritm
- 1 Valet av **dekomposition** bestämmer vilka val som finns senare

Exempel på beroendegraf



Kornighet och parallellitet

- 1 Antalet tasks från dekompositionen och deras storlek bestämmer **kornigheten**
 - **Finkornig** vs **grovkornig**
- 1 Grad av **parallellitet** – "level of concurrency"
 - **Maximum**: största antalet tasks som kan exekveras samtidigt
 - **Genomsnittlig**: genomsnittliga antalet tasks som kan exekveras samtidigt i hela exekveringen
 - Maximum och genomsnittlig grad av parallellitet beror av kornigheten (ökar normalt med kornigheten)

Kornighet och parallellitet

- 1 P = "Den kritiska vägen" för en beroendegraf – längsta vägen mellan något par av start- och målnod.
- 1 L = "Längden av den kritiska vägen", summan av allt arbete längs kritiska vägen
- 1 W = Totala arbetet i parallella algoritmen
- 1 A = W / L, ger **genomsnittliga graden av parallellitet**
- 1 => **Kort kritisk väg ger högre grad av parallellitet**
- 1 Kornigheten kan inte/bör inte ökas hur mycket som helst
 - Problemet kan ha **inbyggda gränser** för kornigheten
 - För hög finkornighet kan ge **andra prestandaproblem**
 - 1 Dåligt cache-utnyttjande
 - 1 För mycket interaktion mellan uppgifterna

Robert Granat

Design and Analysis of Algorithms för Paralleldatorsystem

Interaktion mellan uppgifter

- 1 Normalt förekommer någon form av **interaktion** mellan deluppgifter i alla parallella program
- 1 Kan även ske mellan uppgifter som verkar helt oberoende i programmets beroendegraf
- 1 Interaktionsmönstret kan beskrivas i en **uppgift-interaktions-graf** (noder=tasks, kanter=interaktion)

Robert Granat

Design and Analysis of Algorithms för Paralleldatorsystem

Processer och mappning

- 1 **Dekompositionen ger deluppgifter** som skall exekveras på fysiska processorer
- 1 **Mappning** är den mekanism som **tilldelar uppgifter** till olika processer
- 1 **Process** = mera abstrakt begrepp på **enhet som exekverar kod och använder data hörande till speciell task** inom ramen för den parallella exekveringen
- 1 Tillåter **hierarkisk mappning** av tasks inom flera olika programmeringsparadigmer samtidigt
 - Ex: Message passing mellan noder i paralleldator, där varje nod är en gemensamt-minne maskin med >1 CPU:er (t ex på sarek där varje node är en 2-CPU:ers NUMA)

Robert Granat

Design and Analysis of Algorithms för Paralleldatorsystem

Processer och mappning

- 1 Normalt är **#processer = #processorer**
- 1 En **bra mappning**
 - **Maximera graden av parallellitet**
 - **Minimera totala tiden för beräkningsarbetet** i det parallella programmet
 - **Minimera interaktion mellan processer** i den parallella exekveringen
- 1 En bra mappning uppfyller oftast inte alla dessa krav p g a konflikter, t ex mellan parallellitet och interaktion

Robert Granat

Design and Analysis of Algorithms för Paralleldatorsystem

Rekursiv dekomposition

- 1 Att lösa problemet innebär att **problemet delas** i flera, mindre, problem av **samma sort** (*divide*)
- 1 Lösningen till de mindre problemen måste **kombineras** för att få lösningen till det stora problemet (*conquer*)
- 1 **Gör många algoritmer enkla** att uttrycka
- 1 Varning! Om conquer-steget är stort, kan overheaden bli stor
- 1 Exempel på rekursion: beräkning av explicita matrisinverser av övertriangulära matriser

$$A^{-1} = \begin{bmatrix} A_{11}^{-1} & -A_{11}^{-1}A_{12}A_{22}^{-1} \\ 0 & A_{22}^{-1} \end{bmatrix}$$

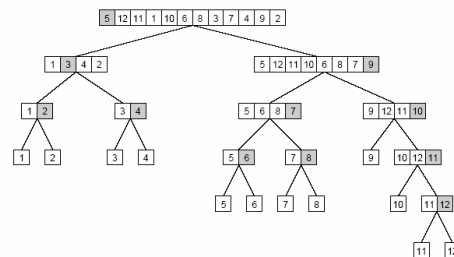
- 1 Beräkningsintensivt conquersteg, dock bibehållen komplexitet
- 1 Lars Karlssons exjobb: <http://www.cs.umu.se/~larsk/>

Robert Granat

Design and Analysis of Algorithms för Paralleldatorsystem

Rekursiv dekomposition

- 1 Ytterligare exempel på rekursiv dekomposition - med litet conquer-steg: quicksort



Robert Granat

Design and Analysis of Algorithms för Paralleldatorsystem

Datadekomposition

I de allra flesta parallella algoritmer är det den **stora datamängden** som är det signifikanta. Två huvudsteg:

1. Datat som beräkningarna skall utföras på **partitioneras**
2. Datapartitioneringen **definierar** en partition av **beräkningsarbetet**

Olika varianter:

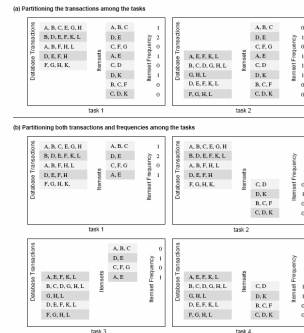
1. Partitionera **utdata**
 - **Oberoende utdata?** Ex: matrismultiplikation $C = A * B$
1. Partitionera **indata**
 - **Hur beräkna utdata?** m h a delresultat från partitionerat indata?
1. Partitionera på **både in- och utdata**
1. Partitionera på **delresultat**
 - Uppdelning av mellanliggande delberäkningar

Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Datadekomposition

Exempel på partition av indata/indata och utdata: **beräkning av frekvenser av grupper av transaktioner i transaktions-databas**

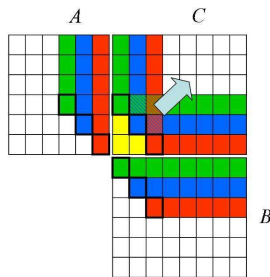


Robert Granat

Design and Analysis of Algorithms for Parallel Computers

”Ägaren beräknar”-regeln

1. Varje datadekomposition av indata och/eller utdata kallas också för **”ägaren beräknar”-regeln**.
1. Iden är att **den som håller en viss del av datat också ansvarar för beräkningarna** som hör till den delen av datat
1. Exempel: vågfronts-algoritmer för $AX - XB = C$
 - Ägaren av C_{ij} ansvarar för beräkningen av X_{ij}
 - C_{ij} skrivs över med X_{ij}

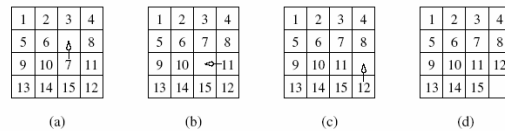


Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Undersökande dekomposition

1. Används vid dekomposition av problem vars lösning beräknas genom **sökning i en lösningsrymd**
1. Grafproblem, spelsökning
1. Dela upp problemområdet i delar vars alla resultat inte **”behövs”**. Sökningen kan **termineras** när någon hittat en (tillräckligt bra) lösning

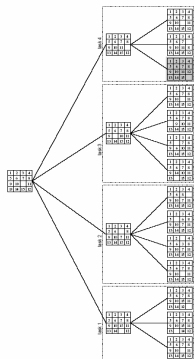


Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Undersökande dekomposition

1. Osäkert hur mycket bättre en parallell algoritm blir: för **mycket jobb kan utföras i onödan** i jämförelse med seriella varianter
1. **Speed-up** närmar sig p endast om man tar **medelvärdet av alla testfall**.

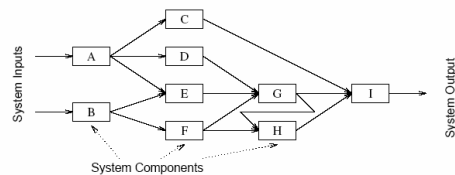


Robert Granat

eddatarsystem

Spekulativ dekomposition

1. Börja utföra beräkningar trots att **alla indata inte är kända**. Ex: börja evaluera alla alternativ i en branch (if, switch) innan villkoret för vägvalet är färdigberäknat
1. Endast lämpligt när indata kan anta några få värden.
1. **Garanterat overhead**.



Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Spekulativ dekomposition

- 1 Andra varianter: evaluera bara de **alternativ** i branchen som verkar **mest troliga**
- 1 Speedup kan bli signifikant om det finns **flera nivåer** av spekulativ dekomposition
- 1 Dock $< p$ eftersom det **alltid** utförs **onödigt arbete**
- 1 Skillnaden mellan undersökande och spekulativ dekomposition:
 - I **undersökande** är **utdata** från multipla uppgifter från en branch **okänd**
 - I **spekulativ** är **indata** till en branch som leder till multipla uppgifter **okänd**

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Hybrid-dekompositioner

- 1 Olika **dekompositionstekniker** kan **kombineras**
- 1 Exempel: dekomposition av matrisberäkning på en parallellmaskin med SMP-noder
 - Datadekomposition av input och/eller output mellan noderna
 - Rekursiv dekomposition av arbetet mellan processorerna inom respektive SMP-nod

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Uppgifter

- 1 Dekompositionen leder till att olika **oberoende** uppgifter (tasks) kan **identifieras**
 - Oberoendet kan ändå innebära **viss interaktion**
- 1 Uppgifterna skall nu **tilldelas de olika processerna**
- 1 Hur skapas uppgifter? **Statiskt** eller **dynamiskt**?
 - Statiskt: **alla uppgifter kända** innan exekveringen startar
 - Dynamiskt: **alla uppgifter inte kända** innan exekveringen startar
 - Dynamiskt skapade uppgifter kräver mer omsorg om lastbalansen och skapar interaktion mellan processer
 - Tilldelning görs vanligen statiskt för statiskt genererade uppgifter och dynamiskt för dynamiskt genererade uppgifter

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Uppgifter

- 1 (Beräknings-) **storlek** på uppgifter?
 - Mängden arbete som krävs för att slutföra uppgiften
 - Uniform/icke uniform – kan påverka lastbalans
- 1 **Vet vi storleken** på uppgifter?
 - Kan användas vid mappning på processerna
- 1 **Storlek på datat** som hör till uppgiften?

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Interaktion mellan uppgifter

Ex: **Kommunikation** eller **hantering av gemensamt minne**

- 1 **Statisk** – interaktioner och tidpunkter kända a priori
- 1 **Dynamisk** – interaktioner och tidpunkter icke kända a priori
- 1 **Reguljär** – interaktioner följer ett givet mönster
- 1 **Irreguljär** – inget givet mönster

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Interaktion mellan uppgifter

- 1 **Enbart läsning** av gemensamt data (read-only)
- 1 **Läsning och skrivning** av gemensamt data (read-write)
- 1 **En-vägs interaktion** initieras och slutförs av en task utan att någon annan blir involverad eller avbruten
 - Kan hanteras av "gemensamt minne"-paradigmer
- 1 **Två-vägs** – "producent och konsument"-scenario
 - Givna modellen för "distribuerat minne"-paradigmer.
 - Förekommer också i "gemensamt minne".

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Mappningstekniker för lastbalansering

- Givet en mängd uppgifter, hur placerar vi dessa på processer för att minimera overhead?
 - Minimera kommunikation (interaktion, synkronisering)
 - Minimera väntetid ("idle")
 - Dessa två mål är ofta i konflikt – finn acceptabel kompromiss
- Förenklat: mappningstekniker statiska eller dynamiska

Robert Granat

Design och Analys av Algoritmer för Paralleldatorsystem

Statisk mappning

- Statisk mappning distribuerar uppgifterna mellan processerna innan exekveringen
- Statisk mappning ofta kombinerat med datadekomposition
 - Block distribution, högre dimension ger generellt högre parallellitet
 - En-dimensionell. Ex: kolumnblocksmappning av 2-dim array
 - Fler-dimensionell. Ex: mappning av 2-dim array i både rad och kolumnblock
 - Block-cyklisk distribution, många fler block än processorer
 - Bra lastbalans
 - Lämplig då olika delar av datat genererar olika mycket arbete, Ex: LU
 - 2D används i ScaLAPACK
 - Cyklisk distribution
 - Rad eller kolumnvis
 - Perfekt lastbalans men bristande lokalitet kan ge sämre prestanda

Robert Granat

Design och Analys av Algoritmer för Paralleldatorsystem

Statisk mappning

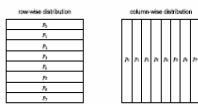


Figure 3.24 Examples of one-dimensional partitioning of an array among eight processes.

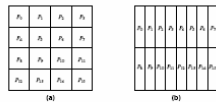


Figure 3.25 Examples of two-dimensional distributions of an array (a) on a 4×4 process grid, and (b) on a 2×8 process grid.

Robert Granat

Design och Analys av Algoritmer för Paralleldatorsystem

Statisk mappning

- Statisk mappning kan också kombineras med slumpmässig blockdistribution
 - Många fler block än processer
 - Blocken delas ut slumpmässigt
 - Kan vara bättre än t ex block-cyklisk vid glesa matriser

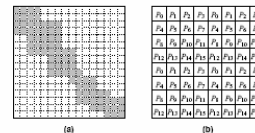


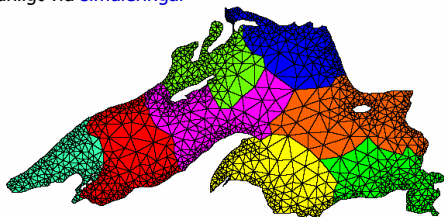
Figure 3.31 Using the block-cyclic distribution shown in (b) to distribute the computations performed in array (a) will lead to load imbalances.

Robert Granat

Design och Analys av Algoritmer för Paralleldatorsystem

Statisk beroendegraf-partitionering

- Dela in datat i delar så att kontaktytorna (=kommunikation) blir så liten som möjligt
- Kontaktytorna bestäms till exempel av en gles matris
- Vanligt vid simuleringar

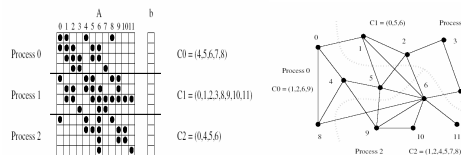


Robert Granat

Design och Analys av Algoritmer för Paralleldatorsystem

Statisk uppgiftspartitionering

- Dela in uppgiftsgrafen i delar så att kontaktytorna (=kommunikation) blir så liten som möjligt
- Ex: gles matris-vektor multiplikation

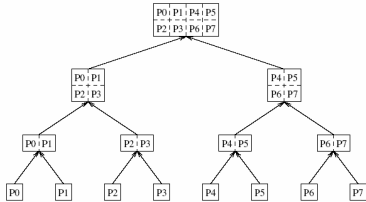


Robert Granat

Design och Analys av Algoritmer för Paralleldatorsystem

Statisk hierarkisk mappning

- Exempel: använd **rekursiv mappning** för att bygga uppgiftsgraf, och **datapartionering** för att ytterligare partitionera de "feta" noderna.



Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Dynamisk mappning

- Nödvändigt då statisk mappning ger obalanser i arbetsbördan mellan olika processer
- Annat namn: **dynamisk lastbalansering**
 - Centraliserad
 - Master-Slave, centraliserad arbetspool
 - Skalar inte bra: masterprocessen blir flaskhals
 - "Chunk scheduling" kan lätta på trycket
 - Distribuerad
 - Vilka par av processer skall utbyta arbete?
 - Vem tar initiativet? Sändaren eller mottagaren?
 - Hur mycket arbete skickas i varje sändning?
 - När skall arbete utbytas? När det är slut? När man vill bli av med lite?
- Kan implementeras i **de flesta paradigmer**

Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Metoder för att dölja interaktion

- Maximera datalokalitet
 - Minimera totala mängden data som utbytes
 - Välj lämplig dekomposition och mappning
 - Ex: distribution av högre dimension oftast bättre
 - Minimera frekvensen av datautbyte
 - Omstrukturering av parallella algoritmen kan vara nödvändig
 - Kommunicera större delar data vid färre tillfällen
- Minimera flaskhalsar
 - Trängsel (contention)
 - Kommunikation
 - Omstrukturering av parallella algoritmen så att beräkningar utföres på klart avgränsade delar

Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Metoder för att dölja interaktion

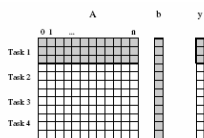
- Överlappa beräkningar med kommunikation
 - Mycket vanligt vid numeriska beräkningar
 - Nackdel: minskar kornigheten på uppgifter
 - Kräver stöd av underliggande mjukvara och/eller hårdvara
 - Distribuerat minne – räkna och kommunicera samtidigt
 - Gemensamt minne – data prefetching
- Använd kollektiva kommunikationsoperationer
 - Broadcast, scatter, reduce, gather, all-to-all
 - Ofta högoptimerade implementationer, t ex *recursive doubling*
- Överlappa interaktioner med andra interaktioner

Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Metoder för att dölja interaktion

- Replikera data
 - Kräver extra minne.
 - Ex: matris-vektor $y = Ab$, alla processer har hela vektorn b



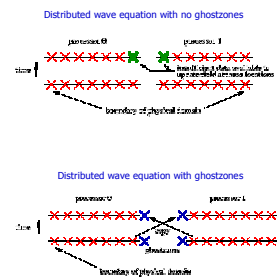
- Utför extra beräkningar (redundant computing)
 - Ex: **spökceller** i beräkningsnät för tidsberoende PDE,
 - Behöver bara kommunicera varannat tidssteg

Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Metoder för att dölja interaktion, exempel

- För att uppdatera värden på inre gränser behövs värden från andra processorer
- Man kan utbyta $O(n)$ datapunkter i varje tidssteg...
- Man kan också välja att utbyta $O(2n)$ datapunkter i vart annat tidssteg och beräkna de $O(n)$ värden som behövs för nästa steg på båda sidor om gränsen
- Extra datat sparas i "spökzoner"
 - Sparar kommunikation – gör dock extra arbete (ofta försurnbart vid grovkornig dekomposition)



Robert Granat

Design and Analysis of Algorithms for Parallel Computers

Modeller för parallella algoritmer

- 1 Data-parallel modell
 - Identiska operationer utföres parallellt på stora datamängder
 - Data-partitionering, statisk mappning, reguljära interaktioner
 - Ex: applicera en beräkningsstencil för PDE över ett diskretiserat område
 - Stora problem kan lösas effektivt
- 1 Uppgifts-parallellism-modell
 - Partitionering av beroendegraf
 - Statisk/dynamisk partitionering och mappning
 - Ex: kompilera oberoende subrutiner, anropa oberoende subrutiner
- 1 Arbetspoolsmodell
 - Dynamisk mappning för god lastbalans
 - Centraliserad eller decentraliserad (se tidigare)
 - Ex: lösa spel- eller grafproblem genom sökning i ett Lösningsrum

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Modeller för parallella algoritmer

- 1 Master-slave
 - En process genererar arbete och distribuerar till arbetare
 - Uppgifterna kan genereras och fördelas statiskt på förhand om möjligt
 - Mastern kan bli flaskhals
 - Ex: distribution av oberoende iterationer i en loop på SMP-maskin
- 1 Pipeline
 - Ström av data passerar genom olika processer
 - Olika processer gör olika saker med datat parallellt
 - Kedja av producent-konsument
 - Ex: Parallell LU-faktorisering
 - Se även fallstudien
- 1 Hybridmodeller
 - Olika modeller appliceras hierarkiskt eller sekventiellt på ett problem

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Fallstudie: MPEG-2 till MPEG-1 transcoding

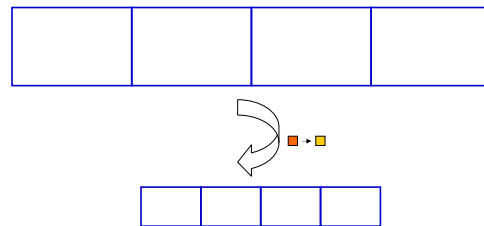
Följande ska utföras:

- 1 En MPEG-2-ström ska omvandlas till en MPEG-1-ström
- 1 Viss korrigering av färgerna ska ske
- 1 MPEG-1-upplösningen är halva MPEG-2-upplösningen
 - Från full TV-upplösning till "skräp-video"
- 1 Både in och ut-format är "long-GOP" (group of pictures)
 - Förenklat:
 - Var 12:e ruta är en referensruta (I)
 - Resten byggs upp med differens från föregående I-frame (P-frames, betydligt mindre än I-frames, P = prediction)
 - I verkliga livet finns också B-frames: differenser som pekar bakåt till närmaste I eller P eller framåt mot närmaste I eller P, B = bi-directional. B-frames möjliggör mera komprimering.
- 1 Typisk lagring av strömmen: IBBPBBPBBPBBPBBPBBPBBP osv...

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Fallstudie: MPEG-2 till MPEG-1 transcoding



Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Fallstudie: MPEG-2 till MPEG-1 transcoding

Dekomposition:

- 1 GOP
- 1 Ruta
- 1 Block
- 1 Hierarkisk
- 1 Uppgift
- 1 ...?

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Fallstudie: MPEG-2 till MPEG-1 transcoding

- 1 Förslag till lösning:
 - Pipeline: alla bilderna strömmar igenom pipelinen där olika processer ansvarar för olika steg
 - 1 Uppackning – bygg bilden av resp I-,P- och B-frames
 - 1 Omskalning
 - 1 Färgkorrigering
 - 1 +några andra operationer (ex. konvertera/mixa ljud)
 - 1 Nedpackning – ta isär bilden till resp I-,P- och B-frames
 - Eventuellt går flera processer in i ett av stegen om det är för kostnadskrävande i jämförelse med de andra stegen
 - 1 Rutorna delas upp blockvis i rektanglar (1-dim datapartitionering)
 - Task-parallellism
 - Lämplig för SMP eller multi(dual/quad)-core maskiner
 - 1 Stöd för både tunga processer (fork) och lättviktstrådar (threads)

Robert Granat

Design och Analys av Algoritmer för Parallellsystem

Sammanfattning – check box

<i>Dekomposition</i>	Rekursion	Data	Undersökande	Spekulativ
<i>Uppgifter</i>	Hur skapas	Storlek	Kunskap om storlek	Storlek på data
<i>Interaktion m. uppgifter</i>	Statisk/ dynamisk	Reguljär/ irreguljär	Läsning/ skrivning	Sändning/ utbyte
<i>Mappning</i>	Statisk:	Array/block/ cyklisk/slump	Graf	Uppgift/ hierarkisk
	Dynamisk:	Centraliserad	Decentraliserad	
<i>Hantera overhead</i>	Maximera datalokalitet	Minimera flaskhalsar	Replikera data/beräkn.	Överlappa/ gruppkommun.
<i>Modell</i>	Dataparallell	Uppgiftsgraf	Arbetspool/ master-slave	Pipeline