Sanders: Parallel Algorithms January 17, 2022



Parallel Algorithms

Peter Sanders Institute of Theoretical Informatics







Why Parallel Processing

increase speed: p computers jointly working on a problem solve it up to p times faster. But, too many cooks spoil the broth. \rightsquigarrow careful coordination of processors

save energy: Two processors with half the clock frequency need less power than one processor at full speed. (power \approx voltage \cdot clock frequency)

expand available memory by using aggregate memory of many processors

less communication: when data is distributed it can also be (pre)processed in a distributed way.

Sanders: Parallel Algorithms January 17, 2022



3

Subject of the Lecture

Basic methods for parallel problem solving

- Parallelization of basic sequential techniques: sorting, data structures, graph algorithms, ...
 - Basic communication patterns
- load balancing
- Emphasis on provable performance guarantees
- But applicability always in view

Overview

- Models, simple examples
- Matrix multiplikation
- Broadcasting
- Sorting
- General data exchange
 - Load balanciong I, II, III
- \Box List ranking (conversion list \rightarrow array)
- Hashing, priority queues
- ☐ Simple graph algorithms



Literature

Script (in German)

+



Figures from the book marked as [Book].





6

More Literature

[Kumar, Grama, Gupta und Karypis],

Introduction to Parallel Computing. Design and Analysis of Algorithms,

Benjamin/Cummings, 1994. mostly practical, programming oriented

[Leighton], Introduction to Parallel Algorithms and Architectures,

Morgan Kaufmann, 1992.

Theoretical algorithms on concrete interconnection networks

[JáJá], <u>An Introduction to Parallel Algorithms</u>, Addison Wesley, 1992. PRAM

[Sanders, Worsch],

Parallele Programmierung mit MPI – ein Praktikum, Logos, 1997.

Parallel Computing at ITI Sanders

- Massively parallel sorting,
- Massively parallel graph algorithms,
- Fault tolerance,
- Shared-memory data structures,
- ☐ (Hyper)graph partitioning,



Michael Axtmann Sebastian Lamm Lukas Hübner Tobias Maier

Tobias Heuer & Daniel Seemaier Lorenz Hübschle-Schneider Markus Iser and Dominik Schreiber Daniel Funke

- Comm. eff. alg.,
- \Box SAT-solving and planning,
- Geometric algorithms,

Role in the CS curriculum

- Elective or "Mastervorzug" Bachelor!
- Vertiefungsfach
 - Algorithmentechnik
 - Parallelverarbeitung
- Studienprofil daten-intensives Rechnen



Sanders: Parallel Algorithms January 17, 2022

Related Courses



9

Parallel programming: Tichy, Karl, Streit

Modelle der Parallelverarbeitung: viel theoretischer,

Komplexitätstheorie,...

Worsch

Algorithmen in Zellularautomaten: spezieller, radikaler, theoretischer Worsch

Rechnerarchitektur: Karl

GPUs: Dachsbacher

+ other algorithms lectures

Sanders: Parallel Algorithms January 17, 2022

RAM/von Neumann Model

Analysis: count machine instructions — load, store, arithmetics, branch,...

simple

very successful









Algorithmenanalyse:

- Count cycles: T(I), for given problem instance I.
- Worst case depending on problem size: $T(n) = \max_{|I|=n} T(I)$

$$\Box \text{ Average case: } T_{\text{avg}}(n) = \frac{\sum_{|I|=n} T(I)}{|\{I : |I|=n\}|}$$

Example: Quicksort has average case complexity $\mathcal{O}(n \log n)$

 \Box Randomized Algorithms: T(n) (worst case) is a random variable.

We are interested, e.g. in its expectation (later more).

Don't mix up with average case.

Example: Quicksort with random pivot selection has expected worst case cost $\mathbb{E}[T(n)] = \mathcal{O}(n \log n)$



Algorithm Analysis: More Conventions

- $\square \mathcal{O}(\cdot)$ gets rid of cumbersome constants
- Secondary goal: memory

Execution time can depend on several parameters:

Example: An efficient variant of Dijkstra's Algorithm for shortest paths needs time $\mathscr{O}(m + n \log n)$ where *n* is the number of nodes and *m* is the number of edges. (It always has to be clear what parameters mean.)



A Simple Parallele Model: PRAM

Idea: change RAM as little as possible.

- □ *p* PEs (**P**rocessing **E**lements); numbered 1..p (or 0..p-1). Every PE knows *p*.
 - One machine instruction per clock cycle and PE synchronous

Shared global memory





Access Conflicts?

EREW: Exclusive Read Exclusive Write. Concurrent access is forbidden

CREW: Concurrent Read Exclusive Write. Concurrent read OK. Example: One PE writes, others read = "Broadcast"

CRCW: Concurrent Read Concurrent Write. avoid chaos:

common: All writers have to agree

Example: OR in constant time (AND?)

arbitrary: someone succeeds

priority: writer with smallest ID succeeds

combine: All values are combined, e.g., sum

 $(\neq random!)$

Example: Global Or

Input in x[1..p]

Initilize memory location Result=0

```
Parallel on PE i = 1..p
```

```
if x[i] then Result := 1
```

Global And

Initilize memory location Result = 1

```
if not x[i] then Result := 0
```



Example: Maximum on Common CRCW PRAM [JáJá Algorithm 2.8]

- Input: A[1..n]// distinct elementsOutput: M[1..n]// M[i] = 1 iff $A[i] = \max_j A[j]$
- forall $(i, j) \in \{1..n\}^2$ dopar $B[i, j] := A[i] \ge A[j]$ forall $i \in \{1..n\}$ dopar $M[i] := \bigwedge_{j=1}^n B[i, j]$ // parallel subroutine

 $\mathscr{O}(1)$ time $\Theta(n^2)$ PEs (!)

16

Sanders: Parallel Algorithms January 17, 2022

i	А	B 1	2	3	4	5 <- j	Karlsruhe Institute of Techno
1	3	*	0	1	0	1	1
2	5	1	*	1	0	1	1
3	2	0	0	*	0	1	1
4	8	1	1	1	*	1	1
5	1	0	0	0	0	*	1
	А	3	5	2	8	1	
 i	–––- A	 В 1	 2	 3	 4	 5 <- j	——— М
 i 1	 А З	 B 1 *	 2 0	 3 1	 4 0	 5 <- j 1	——— М О
 i 1 2	 A 3 5	B 1 * 1	 2 0 *	 3 1 1	 4 0 0	j 1 1	——— М О О
 i 1 2 3	A 3 5 2	B 1 * 1 0	 2 0 * 0	 3 1 1 *	 4 0 0 0	5 <- j 1 1 1	M 0 0 0
 i 1 2 3 4	A 3 5 2 8	B 1 * 1 0 1	2 0 * 0 1	 3 1 1 * 1	4 0 0 0 *	5 <- j 1 1 1 1	 M 0 0 0 1->maxValue=8





18

Describing Parallel Algorithms

- Pascal-like pseudocode
- Explicitly parallel loops [JáJá S. 72]
- Single Program Multiple Data principle. The PE-index is used to break the symmetry. \neq SIMD !



Analysis of Parallel Algorithms

In principle – just another parameter: p. Find execution time T(I, p).

Problem: interpretation.

Work: W = pT(p) is a cost measure. (e.g. Max: $W = \Theta(n^2)$) Span: $T_{\infty} = \inf_p T(p)$ measures parallelizability.

(absolute) Speedup: $S = T_{seq}/T(p)$.

Use the **best known** sequential algorithm.

Relative Speedup $S_{rel} = T(1)/T(p)$ is usually different! (e.g. Maximum: $S = \Theta(n)$, $S_{rel} = \Theta(n^2)$)

Efficiency: E = S/p. Goal: $E \approx 1$ or, at least $E = \Theta(1)$. "Superlinear speedup": E > 1. (possible?). Example maximum: $E = \Theta(1/n)$.

19



20

PRAM vs. Real Parallel Computers

Distributed Memory





(Symmetric) Shared Memory



Problems



Asynchronous ~~> design, analysis, implementation, and debugging are much more difficult than for PRAM

Contention for memory module / cache line

Example:

The $\Theta(1)$ PRAM Algorithm for global OR becomes $\Theta(p)$.

local/cache-memory is (much) faster than global memory

- The network gets more complex with growing *p* while latencies grow
- Contention in the network
- Maximum local memory consumption more important than total memory consumption



23

Realistic Shared-Memory Models

asynchronous

□ aCRQW: asynchronous concurrent read queued write. When *x* PEs contend for the same memory cell, this costs time $\mathscr{O}(x)$.

consistent write operations using atomic operations

memory hierarchies

Why is concurrent read OK?



Atomic Instuctions: Compare-And-Swap

General and widely available:

```
Function CAS(a, expected, desired) : \{0, 1\}BeginTransactionif *a = expected thenelseelseEndTransaction
```

*a:= desired; return 1// success
expected:= *a; return 0// failure



25

Further Operations for Consistent Memory Access:

Fetch-and-add

Hardware transactions

```
Function fetchAndAdd(a, \Delta)
```

```
expected:= *a
```

repeat

```
desired:= expected + \Delta
until CAS(a, expected, desired)
return desired
```



Parallel External Memory





Models with Interconnection Networks



PEs are RAMs

async	hronous	processing

Interaction by message exchange

Important: cost model for data exchange



28

Real Maschines Today





Handling Complex Hierarchies

Proposition: we get quite far with flat Models, in particular for shared memory.

- Design distributed, implement hierarchy adapted
- Shared-memory subroutines on nodes

Explicit ,,Store-and-Forward"

We know the interconnection graph

- $(V = \{1, \dots, p\}, E \subseteq V \times V)$. Variants:
- $V = \{1, \dots, p\} \cup R$ with additional

"dumb" router nodes (perhaps with buffer memory).

- Buses \rightarrow Hyperedges





In each time step each edge can transport up to k' data packets of constant length (usually k' = 1)

On a *k*-port-machine, each node can simultaneously send or receive *k* Packets gleichzeitig senden oder empfangen. k = 1 is called single-ported.



Discussion

- + simple formulation
- low level \Rightarrow "messy algorithms"
- Hardware router allow fast communication whenever an unloaded communication path is found.



Typical Interconnection Networks







hypercube





33

Fully Connected Point-to-Point

 $\Box E = V \times V$, single ported

 $\Box T_{\text{comm}}(m) = \alpha + m\beta$. (*m* = message length in machine words)

- + Realistic handling of message lengths
- + Many interconnection networks approximate fully connected networks \Rightarrow sensible abstraction
- + No overloaded edges \rightarrow OK for hardware router
- + "artificially" increasing lpha, eta
 - \rightarrow OK for "weak" networks
- + Asynchronous model
- A bit of hand waving for real networks

Sanders: Parallel Algorithms January 17, 2022

Fully Connected: Variants



What can PE *I* do in time $T_{\text{comm}}(m) = \alpha + m\beta$? message length *m*.

half duplex: $1 \times \text{send or } 1 \times \text{receive}$ (also called simplex)

Telephone: $1 \times \text{send}$ to PE j and $1 \times \text{receive}$ from PE j

(full)duplex: $1 \times \text{send}$ and $1 \times \text{receive}$.

Arbitrary comunication partners

Effect on running time:

$$T^{\text{duplex}} \leq T^{\text{Telefon}} \leq T^{\text{duplex/2}} \leq 3T^{\text{duplex}}$$

Sanders: Parallel Algorithms January 17, 2022



BSP Bulk Synchronous Parallel

[McColl LNCS Band 1000, S. 46]

Machine described by three parameters: p, L and g.

L: Startup overhead for one collective message exchange – involving all PEs

g: gap $\approx \frac{\text{computation speed}}{\text{communication bandwidth}}$

Superstep: Work locally. Then collective global synchronized data exchange with arbitrary messages.

w: max. local work (clock cycles)

h: max. number of machine words sent or received by a PE (*h*-relation)

Time: w + L + gh

BSP versus Point-to-Point

By naive direct data delivery:

Let $H = \max$ #messages of a PEs.

Then $T \ge \alpha(H + \log p) + h\beta$.

Worst case H = h. Thus $L \ge \alpha \log p$ and $g \ge \alpha$?

By all-to-all and direct data delivery:

Then $T \ge \alpha p + h\beta$. Thus $L \ge \alpha p$ and $g \approx \beta$?

By all-to-all and indirect data delivery:

Then $T = \Omega(\log p(\alpha + h\beta))$. Thus $L = \Omega(\alpha \log p)$ and $g = \Omega(\beta \log p)$?


BSP*



Armin Bäumker and Friedhelm Meyer auf der Heide, ESA 1995.

Redefinition of *h* to # blocks of size *B*, e.g. $B = \Theta(\alpha/\beta)$.

Let M_i be the set of messages sent or received by PE i. Let $h = \max_i \sum_{m \in M_i} \lceil |m| / B \rceil$.

Let g the gap between sending packets of size B.

Then, once more

$$w + L + gh$$

is the time for a super step.

BSP* versus Point-to-Point

With naive direct data delivery:

 $L \approx \alpha \log p$ $g \approx B\beta$



BSP⁺

We augment BSP by collective operations

broadcast

(all-)reduce

prefix-sum

with message lengths h.

Stay tuned for algorithms which justify this.

 $\mathsf{BSP}^*\text{-}\mathsf{algorithms}$ are up to a factor $\Theta(\log p)$ slower than $\mathsf{BSP}^+\text{-}\mathsf{algorithms}.$





MapReduce





MapReduce Example WordCount



MapReduce Discussion

- + Abstracts away difficult issues
 - * parallelization
 - * load balancing
 - * fault tolerance
 - * memory hierarchies
 - * .
- Large overheads
- Limited functionality



$$A \subseteq I$$

1: map

$$B = \bigcup_{a \in A} \mu(a) \subseteq K \times V$$

2: shuffle

$$C = \{(k, X) : k \in K \land$$

3: reduce

$$X = \{x : (k, x) \in B\} \land X \neq \emptyset$$

$$D = \bigcup_{c \in C} \rho(c)$$

MapReduce – MRC Model

[Karloff Suri Vassilvitskii 2010]

A problem is in MRC iff for input of size *n*:

 \Box solvable in $\mathscr{O}(\operatorname{polylog}(n))$

MapReduce steps

 $\ \square \ \mu$ and ρ evaluate in time $\mathscr{O}(\mathrm{poly}(n))$

 $\square \ \mu \text{ and } \rho \text{ use space } \mathscr{O}(n^{1-\varepsilon})$ ("substantially sublinear")

□ overall space for $B \mathscr{O}(n^{2-\varepsilon})$ ("substantially subquadratic")

Roughly: count steps,

very loose constraints on everything else



$$A \subseteq I$$

1: map

$$B = \bigcup_{a \in A} \mu(a) \subseteq K \times V$$

2: shuffle

$$C = \{(k, X) : k \in K \land$$

3: reduce

$$X = \{x : (k, x) \in B \land$$

$$X \neq \emptyset\}$$

$$D = \bigcup_{c \in C} \rho(c)$$

MapReduce – MRC Model



A problem is in MRC iff for input of size *n*:

solvable in $\mathscr{O}(\operatorname{polylog}(n))$

MapReduce steps

 \square μ and ρ evaluate in time $\mathscr{O}(\mathrm{poly}(n))$

 $\square \ \mu \text{ and } \rho \text{ use space } \mathscr{O}(n^{1-\varepsilon})$ ("substantially sublinear")

 \Box overall space for $B \ \mathscr{O} \left(n^{2-\varepsilon} \right)$

("substantially subquadratic")

Roughly: count steps,

very loose constraints on everything else

 n^2 ? n^{42} ? "big" data?

"big" data?

speedup? efficiency?

MapReduce – MRC+ Model

[Sanders IEEE BigData 2020]

- w Total work for μ , ρ
- \hat{w} Maximal work for μ , ρ
- Total data volume for $A \cup B \cup C \cup D$ m
- \hat{m} Maximal object size in $A \cup B \cup C \cup D$,

Theorem: Lower and upper bound for parallel exec $\Theta\left(\frac{w}{p} + \hat{w} + \log p\right) \text{ parallel time,} \\ \Theta\left(\frac{m}{p} + \hat{m} + \log p\right) \text{ bottleneck communication volume.}$



$$A \subseteq I$$
1: map
$$B = \bigcup_{a \in A} \mu(a) \subseteq K \times V$$
2: shuffle
$$C = \{(k, X) : k \in K \land$$
3: reduce
$$X = \{x : (k, x) \in B \land$$

$$X \neq \emptyset\}$$

$$D = \bigcup_{c \in C} \rho(c)$$
cution:



46

Implementing MapReduce on BSP

2 Supersteps:

1. Map local data.

Send $(k, v) \in B$ to PE h(k) (hashing)

2. Receive elements of B.

Build elements of *C*.

Run reducer.

Send all but first result of a reducer to a random PE.



MapReduce on BSP – Example



47



48

MapReduce on BSP – Analysis

Time

$$2L + \mathscr{O}\left(\frac{w}{p} + g\frac{m}{p}\right)$$

if

$$w = \Omega(\hat{w}p\log p)$$
 and $m = \Omega(\hat{m}p\log p)$

Graph- und Circuit Representation of Algorithms

Many computations can be represented as a directed acyclic graph

Input nodes have indegree 0 and a fixed output

Ouput nodes have outdegree 0 and indegree 1

- The indegree is bounded by
 - a small constant.
- Inner nodes compute a function, that can be computed in constant time.



Circuits



- Variant: When a constant number of bits rather than machine words are processed, we get circuits.
- The depth d(S) of the computation-DAG is the number of inner nodes on the longest path from an input to an output \sim time
- One circuit for each input size (specified algorithmically) \Rightarrow circuit family

Example: Associative Operations (=Reduction)

Satz 1. Let \oplus denote an associative operator that can be computed in constant time. Then,

$$\bigoplus_{i < n} x_i := (\cdots ((x_0 \oplus x_1) \oplus x_2) \oplus \cdots \oplus x_{n-1})$$

can be computed in time $\mathcal{O}(\log n)$ on a PRAM and in time $\mathcal{O}(\alpha \log n)$ on a linear array with hardware router

Example: $+, \cdot, \max, \min, \ldots$ (example for non-commutative?)



Proof Outline for
$$n = 2^k$$
 (wlog?)

Induction hypothesis: \exists circuit of depth *k* for $\bigoplus_{i < 2^k} x_i$

k = 0: trivial



PRAM Code PE index $i \in \{0, ..., n-1\}$ active := 1Х for $0 \le k < \lceil \log n \rceil$ do if active then if bit k of i then active := 0else if $i + 2^k < n$ then $x_i \coloneqq x_i \oplus x_{i+2^k}$ //result is in x_0

Careful: much more complicated on a real asynchronous shared-memory machine.

Speedup? Efficiency?

 $\log x$ here always $\log_2 x$

Analysis

n PEs Time $\mathcal{O}(\log n)$ Speedup $\mathcal{O}(n/\log n)$ Efficiency $\mathcal{O}(1/\log n)$



54



Less is More (Brent's Principle)

p PEs

Each PE adds

n/p elements sequentially.

Then parallel sum

for p subsums

Time $T_{\text{seq}}(n/p) + \Theta(\log p)$ Efficiency



$$\frac{T_{\text{seq}}(n)}{p(T_{\text{seq}}(n/p) + \Theta(\log p))} = \frac{1}{1 + \Theta(p\log(p))/n} = 1 - \Theta\left(\frac{p\log p}{n}\right)$$

if $n \gg p\log p$

55



Distributed Memory Machine

PE index $i \in \{0, \dots, n-1\}$ //Input x_i located on PE i

active := 1

 $S := X_i$

for $0 \le k < \lceil \log n \rceil$ do

if active then

if bit k of i **then** sync-send s to PE $i - 2^k$

active := 0 else if $i + 2^k < n$ then receive s' from PE $i + 2^k$

 $s := s \oplus s'$

//result is in s on PE 0



Analysis



fully connected: $\Theta((\alpha + \beta) \log p)$

linear array: $\Theta(p)$: step k needs time 2^k .

linear array with router: $\Theta((\alpha+\beta)\log p),$ since edge congestion is one in every step

 $\mathsf{BSP}\ \Theta((l+g)\log p) = \Omega\left(\log^2 p\right)$

Arbitrary n > p: additional time $T_{seq}(n/p)$



58

Discussion Reductions Operation

- Binary tree yields logarithmic running time
- Useful for most models
- Brent's principle: inefficient algorithms get efficient by using less
 PEs
- Later: reduction of complex objects, e.g., vectors, matrices

Matrix Multiplikation

Given: Matrices $A \in \mathbf{R}^{n \times n}$, $B \in \mathbf{R}^{n \times n}$ with $A = ((a_{ij}))$ und $B = ((b_{ij}))$ **R**: semiring

 $C = ((c_{ij})) = A \cdot B$ well-known:

$$c_{ij} = \sum_{k=1}^{n} a_{ik} \cdot b_{kj}$$

work: $\Theta(n^3)$ arithmetical operations (better algorithms if *R* allows subtraction)



A first PRAM Algorithm

 n^3 PEs

for i := 1 to n dopar for j := 1 to n dopar $c_{ij} := \sum_{k=1}^{n} a_{ik} \cdot b_{kj}$

// *n* PE parallel sum

One PE für each product $c_{ikj} := a_{ik}b_{kj}$ Time $\mathscr{O}(\log n)$ Efficiency $\mathscr{O}(1/\log n)$



Distributed Implementation I

 $p \leq n^2 \; \mathrm{PEs}$

for *i*:= 1 to *n* dopar for *j*:= 1 to *n* dopar $c_{ij} := \sum_{k=1}^{n} a_{ik} \cdot b_{kj}$

Assign n^2/p of the c_{ij} to each PE

- limited scalability

- high communication volume. Time
$$\Omega\left(\beta \frac{n^2}{\sqrt{p}}\right)$$





Distributed Implementation II-1

[Dekel Nassimi Sahni 81, KGGK Section 5.4.4]

Assume $p = N^3$, *n* is a multiple of *N* View *A*, *B*, *C* as $N \times N$ matrices, each element is a $n/N \times n/N$ matrix

for i := 1 to N dopar for j := 1 to N dopar $c_{ij} := \sum_{k=1}^{N} a_{ik} b_{kj}$

One PE for each product $c_{ikj} := a_{ik}b_{kj}$







Distributed Implementation II-2

store a_{ik} in PE (i,k,1)

store b_{kj} in PE (1,k,j)

 $\begin{array}{l} \mathsf{PE}\ (i,k,1) \text{ broadcasts } a_{ik} \text{ to } \mathsf{PEs}\ (i,k,j) \text{ for } j \in \{1..N\} \\ \mathsf{PE}\ (1,k,j) \text{ broadcasts } b_{kj} \text{ to } \mathsf{PEs}\ (i,k,j) \text{ for } i \in \{1..N\} \\ \mathsf{compute}\ c_{ikj} \coloneqq a_{ik} b_{kj} \text{ on } \mathsf{PE}\ (i,k,j) \end{array} \right. \qquad \qquad \textit{II local!}$

 $\mathsf{PEs}\ (i, k, j) \text{ for } k \in \{1..N\} \text{ compute } c_{ij} \coloneqq \sum_{k=1}^{N} c_{ikj} \text{ to } \mathsf{PE}\ (i, 1, j)$





64

Analysis, Fully Connected, etc. store a_{ik} in PE (i,k,1) // free (or cheap) store b_{kj} in PE (1,k,j) // free (or cheap) PE (i,k,1) broadcasts a_{ik} to PEs (i,k,j) for $j \in \{1..N\}$ PE (1,k,j) broadcasts b_{kj} to PEs (i,k,j) for $i \in \{1..N\}$ compute $c_{ikj} := a_{ik}b_{kj}$ on PE (i,k,j) // $T_{seq}(n/N) = \mathcal{O}((n/N)^3)$ PEs (i,k,j) for $k \in \{1..N\}$ compute $c_{ij} := \sum_{k=1}^{N} c_{ikj}$ to PE (i,1,j)

Communication:

 $2T_{\text{broadcast}}((n/N)^2, N) + T_{\text{reduce}}((n/N)^2, N) \approx 3T_{\text{broadcast}}((n/N)^2, N)$ $\stackrel{N=p^{1/3}}{\leadsto} \cdots \mathscr{O}\left(\frac{n^3}{p} + \beta \frac{n^2}{p^{2/3}} + \alpha \log p\right)$



Discussion Matrix Multiplikation

- **PRAM** Alg. is a good starting point
- DNS algorithm saves communication but needs factor $\Theta(\sqrt[3]{p})$ more space than other algorithms
- → good for small matrices (for big ones communication is irrelevant)
 - **Pattern** for dense linear algebra:
 - Local Ops on submatrices + Broadcast + Reduce
 - e.g. matrix-vector-product, solve lin. eq. syst.,...

Broadcast and Reduction

Broadcast: One for all

One PE (e.g. 0) sends message of length n to all PEs



Reduction: One for all

One PE (e.g. 0) receives sum of p messages of length n (vector addition \neq local addition)



Broadcast ~>> **Reduction**

turn around direction of communication

add corresponding parts of arriving and own messages

All the following

broadcast algorithms yield

reduction algorithms

for commutative and associative operations.

Most of them (except Johnsson/Ho and some network embeddings) also work for noncommutative operations.







Modelling Assumptions

fully connected

☐ fullduplex – parallel send and receive

Variants: halfduplex, i.e., send or receive, BSP, embedding into concrete networks



Naive Broadcast [KGGK Section 3.2.1]

Procedure naiveBroadcast(m[1..n])

- PE 0: for i := 1 to p 1 do send m to PE i
- PE i > 0: receive m

Time: $(p-1)(n\beta+lpha)$

nightmare for implementing a scalable algorithm





Binomial Tree Broadcast

 $\begin{array}{ll} \textbf{Procedure binomialTreeBroadcast}(m[1..n]) \\ \textbf{PE index } i \in \{0, \ldots, p-1\} \\ \textit{//Message } m \text{ located on PE 0} \\ \textbf{if } i > 0 \textbf{ then receive } m \\ \textbf{for } k := \min\{\lceil \log n \rceil, \operatorname{trailingZeroes}(i)\} - 1 \textbf{ downto } 0 \textbf{ do} \\ & \text{ send } m \text{ to PE } i + 2^k \qquad // \text{ noop if receiver} \ge p \end{array}$



Analysis

- \Box Time: $\lceil \log p \rceil (n\beta + \alpha)$
- \Box Optimal for n = 1

Embeddable into linear array

 $n \cdot f(p) \rightsquigarrow n + \log p$?





Linear Pipeline

Procedure linearPipelineBroadcast(m[1..n], k)

PE index $i \in \{0,\ldots,p-1\}$

//Message m located on PE 0

l/assume k divides n

define piece j as $m[(j-1)\frac{n}{k}+1..j\frac{n}{k}]$

for j := 1 to k + 1 do





receive piece *j* from PE i - 1 // noop if i = 0 or j = k + 1 and, concurrently,

send piece j - 1 to PE i + 1 // noop if i = p - 1 or j = 1
Analysis

 $\Box \text{ Time } \frac{n}{k}\beta + \alpha \text{ per } \underline{\text{step}}$ $(\neq \text{ Iteration})$

 $\Box p-1$ steps until first packet arrives

Then 1 step per further packet

$$T(n, p, k): \left(\frac{n}{k}\beta + \alpha\right)(p + k - 2))$$

optimal k: $\sqrt{\frac{n(p-2)\beta}{\alpha}}$
 $T^*(n, p): \approx n\beta + p\alpha + 2\sqrt{np\alpha\beta}$







74





Discussion



 \Box But for large *p* extremely large messages needed

 $\alpha p \rightsquigarrow \alpha \log p$?





Procedure binaryTreePipelinedBroadcast(m[1..n], k)//Message *m* located on root, assume *k* divides *n* define piece *j* as $m[(j-1)\frac{n}{k}+1..j\frac{n}{k}]$ for *j* := 1 to *k* do if parent exists then receive piece *j* if left child ℓ exists then send piece *j* to ℓ if right child *r* exists then send piece *j* to *r*





Example



Analysis



$$\Box$$
 time $\frac{n}{k}\beta + \alpha$ per step (\neq iteration)

 $\Box 2j$ steps until first packet reaches level j

 \Box how many levels? $d := \lfloor \log p \rfloor$

 \Box then **3** steps for each further packet

Overall:
$$T(n, p, k) := (2d + 3(k-1))\left(\frac{n}{k}\beta + \alpha\right)$$

optimal
$$k$$
: $\sqrt{\frac{n(2d-3)\beta}{3\alpha}}$

Analysis



 \Box *d*:= $\lfloor \log p \rfloor$ levels

$$\Box \text{ Overall: } T(n, p, k) := (2d + 3(k - 1)) \left(\frac{n}{k}\beta + \alpha\right)$$

$$\Box$$
 optimal $k: \sqrt{\frac{n(2d-3)\beta}{3\alpha}}$

substituted: $T^*(n,p) = 2d\alpha + 3n\beta + \mathcal{O}\left(\sqrt{nd\alpha\beta}\right)$





Fibonacci-Trees





Analysis

$$\Box$$
 Time $\frac{n}{k}\beta + \alpha$ per step (\neq iteration)

 \Box *j* steps until first packet reaches level *j*

How many PEs
$$p_j$$
 with level $0..j$?
 $p_0 = 1, p_1 = 2, p_j = p_{j-2} + p_{j-1} + 1 \rightarrow \text{ask Maple},$
 $rsolve(p(0)=1, p(1)=2, p(i)=p(i-2)+p(i-1)+1, p(i));$
 $p_j \approx \frac{3\sqrt{5}+5}{5(\sqrt{5}-1)} \Phi^j \approx 1.89 \Phi^j$
with $\Phi = \frac{1+\sqrt{5}}{2}$ (golden ratio)

 $\rightsquigarrow d \approx \log_{\Phi} p$ levels

overall:
$$T^*(n,p) = d\alpha + 3n\beta + \mathcal{O}\left(\sqrt{nd\alpha\beta}\right)$$



83

Procedure fullDuplexBinaryTreePipelinedBroadcast(m[1..n], k)//Message *m* located on root, assume *k* divides *n* define piece *j* as $m[(j-1)\frac{n}{k}+1..j\frac{n}{k}]$ for *j* := 1 to *k*+1 do receive piece *j* from parent // noop for root or *j* = *k*+1 and, concurrently, send piece *j* - 1 to right child // noop if no such child or *j* = 1

send piece j to left child

– even step

// noop if no such child or j = k + 1

Analysis

- \Box Time $rac{n}{k}eta+lpha$ per step
- $\exists j$ steps until first packet level j reaches
- $\Box \ d \approx \log_{\Phi} p$ levels

Then 2 steps for each further packet

Overall: $T^*(n,p) = d\alpha + 2n\beta + \mathcal{O}\left(\sqrt{nd\alpha\beta}\right)$











86

Discussion

Karlsruhe Institute of Technology

87

Fibonacci trees are a good compromise for all n, p.

General *p*:

use next larger tree. Then drop a subtree



H-Trees







88

H-Trees





••	•	-•	•	•	•	•
• •	•	•	•	-•	•	•
••	•	•	•	•	•	-•
••	•			-•		-0
• •	•	•	•	•	•	_
	•					_
				-		





Disadvantages of Tree-based Broadcasts

Leaves only receive their data. Otherwise they are unproductive

Inner nodes send more than they receive

 \rightsquigarrow full-duplex communication not fully exploited



Binary-Tree-Broadcasts using two trees A und B at once

Inner nodes of A are

leaves of B

and vice versa

Per double step:

One packet received as a leaf +one packet received and forwarded

i.e., 2 packets sent and received

0 0 0 () () 5 8 10 11 12 13 14 3 0 2 4 7 9 6 as inner node 0



Root Process

for *j* := 1 **to** *k* **step 2 do**

send piece j + 0 along edge labelled 0 send piece j + 1 along edge labelled 1





Other Processes,

Wait for first piece to arrive

 \mathbf{if} it comes from the upper tree over an edge labelled b then

 $\Delta := 2 \cdot$ distance of the node from the bottom in the upper tree

for j := 1 to $k + \Delta$ step 2 do

along *b*-edges: receive piece j and send piece j-2

along 1 - b-edges: receive piece $j + 1 - \Delta$ and send piece j





















Arbitrary Number of Processors



98



Constructing the Trees

Case $p = 2^{h} - 1$: Upper tree + lower tree + root Upper tree: complete binary tree of height h - 1, - right leaf Lower tree: complete binary tree of height h - 1, - left leaf Lower tree \approx Upper tree shifted by one. Inner nodes upper tree = leaves of lower tree. Inner nodes lower tree = leaves of upper tree.



Building Smaller Trees (without root)

invariant : last node has outdegree 1 in tree x**invariant** : last node has outdegree 0 in tree \bar{x}

 $p \rightsquigarrow p-1$:

Remove last node:

right node in x now has degree 0 right node in \overline{x} now has degree 1







Coloring Edges

Consider bipartite graph:

$$B = (\{s_0, \dots, s_{p-1}\} \cup \{r_0, \dots, r_{p-2}\}, E)$$

- s_i : Sender role of PE *i*.
- r_i : Receiver role of PE i.
- $2 \times$ degree 1. all other degree 2.
- \Rightarrow *B* is a path plus even circles.

color edges alternatingly with 0 and 1.







Open Question: Parallel Coloring ?

In Time Polylog(p) using list ranking.

(unfortunately impractical for small inputs)

 \Box Fast explicit calculation color(i, p) without communication ?





Jochen Speck's Solution

//Compute color of edge entering node *i* in the upper tree.

I/*h* is a lower bound on the height of node *i*.

Function in EdgeColor(p, i, h)

if *i* is the root of T_1 **then return** 1 **while** *i* bitand $2^h = 0$ **do** h + + // compute height $i' := \begin{cases} i - 2^h & \text{if } 2^{h+1} \text{ bitand } i = 1 \lor i + 2^h > p \\ i + 2^h & \text{else} \end{cases}$ // compute parent of *i*

return in EdgeColor $(p, i', h) \operatorname{xor} (p/2 \mod 2) \operatorname{xor} [i' > i]$



Analysis

- $\Box 2j$ steps until all PEs in level j are reached
- $\Box \ d = \lceil \log(p+1) \rceil \text{ levels}$

□ Then 2 steps for 2 further packets

$$T(n, p, k) \approx \left(\frac{n}{k}\beta + \alpha\right)(2d + k - 1))$$
, with $d \approx \log p$
optimal k : $\sqrt{\frac{n(2d - 1)\beta}{\alpha}}$

 $T^*(n,p) \approx n\beta + \alpha \cdot 2\log p + \sqrt{2n\log p\alpha\beta}$











Implementation in the Simplex-Model

2 steps duplex \rightsquigarrow 4 steps simplex.

1 PE duplex \rightsquigarrow 1 simplex couple = sender + receiver. $0 \downarrow \downarrow 1$





23-Reduktion

Numbering is an inorder-numbering for both trees !




Another optimal Algorithm

[Johnsson Ho 85: Optimal Broadcasting and Personalized Communication in Hypercube, IEEE Transactions on Computers, vol. 38, no.9, pp. 1249-1268.]

Model: full-duplex limited to a single edge per PE (telephone model) Adaptation half-duplex: everything $\times 2$



Hypercube *H*_d

$$\square p = 2^d \mathsf{PEs}$$

 \Box nodes $V = \{0, 1\}^d$, i.e., write node number in binary

 \Box edges in dimension $i: E_i = \left\{ (u, v) : u \oplus v = 2^i \right\}$

$$\Box E = E_0 \cup \cdots \cup E_{d-1}$$





ESBT-Broadcasting

 \Box In step *i* communication along dimension *i* mod *d*

 \Box Decompose H_d into d Edge-disjoint Spanning Binomial Trees

 $] 0^d$ cyclically distributes packet to roots of the ESBTs

ESBT-roots perform binomial tree broadcasting (except missing smallest subtree 0^d)



Analysis, Telephone model

k packets, k divides n

 $\Box k$ steps until last packet left PE 0

 $\exists d$ steps until it has reached the last leaf

 \Box Overall d + k steps

$$T(n, p, k) = \left(\frac{n}{k}\beta + \alpha\right)(k+d)$$

optimal $k: \sqrt{\frac{nd\beta}{\alpha}}$

 $T^*(n,p) := n\beta + d\alpha + \sqrt{nd\alpha\beta}$





Discussion





Reality Check

Libraries (z.B. MPI) often do not have a pipelined implementation of collective operations ~> your own broadcast may be significantly faster than a library routine

Choosing k is more complicated: only piece-wise linear cost function for point-to-point communication, rounding

- Hypercube gets slow when communication latencies have large variance
- perhaps modify Fibonacci-tree etc. in case of asynchronous communication (Sender finishes before receiver). Data should reach all leaves about at the same time.



Broadcast for Library Authors

- \Box **One** Implementation? \rightsquigarrow 23-Broadcast
- Few, simple Variants? {binomial tree, 23-Broadcast} or {binomial tree, 23-Broadcast, linear pipeline}



Beyond Broadcast

Pipelining is important technique for handling large data sets

- hyper-cube algorithms are often elegant and efficient. (and often simpler than ESBT)
 - Parameter tuning (e.g. for. k) is often important.

Sorting

- Fast inefficient ranking
- Quicksort
- □ Sample Sort
- Multiway Mergesort
- ☐ Selection
- More on sorting



Fast Inefficient Ranking



Input: A[1..n]Output: M[1..m] // distinct elements
// M[i] =rank of A[i]

forall
$$(i, j) \in \{1..n\}^2$$
 dopar $B[i, j] := A[i] \ge A[j]$
forall $i \in \{1..n\}$ dopar
 $M[i] := \sum_{j=1}^n B[i, j]$ // parallel subroutine

Running time: $\approx T_{\text{broadcast}}(1) + T_{\text{reduce}}(1) = \mathscr{O}(\alpha \log p)$







Sorting Larger Data Sets



Perhaps more general initial state

Output is globally sorted

$$s_{0,0} \le \dots \le s_{0,n_1-1} \le \dots \le s_{p-1,0} \le \dots \le s_{p-1,n_{p-1}-1}$$

Comparison based model

$$\Box T_{\text{seq}} = T_{\text{compr}} n \log n + \mathcal{O}(n)$$

Careful: different notation in Skript $n \leftrightarrow n/p$



Back to Fast Ranking

//Assume $p = a \times b$ PEs, PE index is (i, j)

Procedure matrix Rank(s)

sort(s) r:= all-gather-by-rows(s, merge) c:= all-gather-by-cols(s, merge) ranks:= $\langle | \{x \in c : x \leq y\} | : y \in r \rangle$ reduce-by-rows(ranks)



// merge

Time

$$\mathscr{O}\left(\alpha\log p + \beta\frac{n}{\sqrt{p}} + \frac{n}{p}\log\frac{n}{p}\right) \quad . \tag{1}$$





Example





123

row all-gather-merge





b

m

a

e

k



row all-gather-merge col all-gather-merge



124





dghl e 0123 h k	dghl a 1223 g i	dghl C 2222 m	dghl b 1113 J 1
abem <mark>e</mark>	abema	abem C	abemb
0013 ^h	1113 g	0023 d	0113J
k	i	m	1
cijk <mark>e</mark>	cijk <mark>a</mark>	cijkc	cijkb
0223 h	1333 g	1222d	1122 J
k	i	m	1





dghl e 0123 h k	dghl a 1223 g i	dghl <mark>C</mark> 2222 m	dghl b 1113 J 1	d 4	д С	h 7	1 11
abem <mark>e</mark> 0013 <mark>k</mark>	abem ^a 1113 ⁹ i	abem C 0023 m	abemb 0113J 1	a 1	b 2	e 5	m 12
cijk <mark>e</mark> 0223 <mark>k</mark>	cijk <mark>a</mark> 1333 g i	cijkc 1222d m	cijkb 1122j 1	C 3	: i 8	ј 9	k 10



More Accurate Analysis

local sorting:
$$\frac{n}{p} \log \frac{n}{p} T_{\text{compr}}$$

2× **all-gather:** $2\left(\beta n/\sqrt{p} + \frac{1}{2}\alpha \log p\right)$

local ranking: $2T_{compr}n/\sqrt{p}$

reduce EDSBT-Algorithm:

$$\beta n/\sqrt{p} + \frac{1}{2}\alpha \log p + \sqrt{\alpha\beta n/\sqrt{p}\frac{1}{2}\log p}$$

Overall:

$$\frac{3}{2}\log p\alpha + 3\beta n/\sqrt{p} + \sqrt{\alpha\beta 0.5n/\sqrt{p}\log p} + \frac{n}{p}\log \frac{n}{p}T_{\text{compr}}$$



Numerical Example:

$$p = 1024$$
, $\alpha = 10^{-5}$ s, $\beta = 10^{-8}$ s, $T_{compr} = 10^{-8}$ s, $n/p = 32$.

$$\frac{3}{2}\log p\alpha + 3n\sqrt{p}\beta + \sqrt{0.5n\sqrt{p}\log p\alpha\beta} + n\log nT_{\rm compr}$$

Time ≈ 0.200 ms.

In comparison: efficient Gather+seq. sort:

 $2 \cdot 32000 \cdot 10^{-8} + 10 \cdot 10^{-5} + 32000 \cdot 15 \cdot 10^{-8} \approx 5.6$ ms

even larger difference with naive gather





129

Quicksort

Sequential

Procedure qSort(d[], n)

if n = 1 then return

select a pivot v

reorder the elements in d such that

$$d_0 \cdots d_{k-1} \le v < d_k \cdots d_{n-1}$$

qSort($[d_0, \dots, d_{k-1}], k$)
qSort($[d_{k+1}, \dots, d_{n-1}], n-k-1$)



Parallelization for Beginners

Parallelization of recursive calls.

 $T_{\rm par} = \Omega(n)$

Very limited speedup

Bad for distributed memory





Parallelization for Theoreticians

For simplicity: n = p.

Idea: Also parallelize partitioning

- 1. One PE provides the pivot (e.g. random choice).
- 2. Broadcast
- 3. Local comparison
- 4. Enumerate "small" elements (prefix-sum)
- 5. Redistribute data
- 6. Split PEs
- 7. Parallel recursion



Parallelization for Theoreticians

//Let $i \in 0..p - 1$ and p denote the 'local' PE index and partition size **Procedure** theoQSort(d, i, p)

if p = 1 then return

j:= random element from 0..p - 1// same value in entire partition v := d@j// broadcast pivot f := d < v $j := \sum_{k=0}^{i} f@k$ // prefix sum p' := j@(p-1)// broadcast if f then send d to PE jelse send d to PE p'+i-j // $i-j = \sum_{k=0}^{i} d@k > v$ receive dif i < p' then join left partition; qsort(d, i, p')

else join right partition; qsort(d, i - p', p - p')



Example

pivot v = 44





```
int pQuickSort(int item, MPI_Comm comm)
{ int iP, nP, small, allSmall, pivot;
   MPI_Comm newComm; MPI_Status status;
   MPI_Comm_rank(comm, &iP); MPI_Comm_size(comm, &nP);
```

```
if (nP == 1) { return item; }
else {
  pivot = getPivot(item, comm, nP);
  count(item < pivot, &small, &allSmall, comm, nP);</pre>
  if (item < pivot) {
    MPI_Bsend(&item,1,MPI_INT, small - 1 ,8,comm);
  } else {
    MPI_Bsend(&item, 1, MPI_INT, allSmall+iP-small, 8, comm);
  }
  MPI Recv(&item,1,MPI INT,MPI ANY SOURCE,8,comm,&status);
  MPI_Comm_split(comm, iP < allSmall, 0, &newComm);</pre>
  return pQuickSort(item, newComm); } }
```



```
/* determine a pivot */
int getPivot(int item, MPI_Comm comm, int nP)
{ int pivot = item;
    int pivotPE = globalRandInt(nP);/* from random PE */
    /* overwrite pivot by that one from pivotPE */
    MPI_Bcast(&pivot, 1, MPI_INT, pivotPE, comm);
    return pivot;
}
```

/* determine prefix-sum and overall sum over value */
void

count(int value, int *sum, int *allSum, MPI_Comm comm, int nP)

{ MPI_Scan(&value, sum, 1, MPI_INT, MPI_SUM, comm); *allSum = *sum; MPI_Bcast(allSum, 1, MPI_INT, nP - 1, comm); }



Analysis

- Per recursion level:
 - $2 \times$ Broadcast
 - $1 \times$ Prefix sum (\rightarrow later)
- \rightsquigarrow Time $\mathscr{O}(\alpha \log p)$

 $\Box \text{ Expected recursion depth: } \mathcal{O}(\log p)$ $(\rightarrow \text{ lecture randomized algorithms})$

Expected overall time: $\mathcal{O}(\alpha \log^2 p)$



Generalization for $n \gg p$ **by Standard Procedure?**

Each PE has, in general, "large" und "small" Elements.

 \Box These numbers are not multiples of n/p

Prefix sums remain useful

$$\Box$$
 On PRAMs we get a $\mathscr{O}\left(rac{n\log n}{p} + \log^2 p\right)$ algorithm

For distributed memory its bad that many elements get moved $\Omega(\log p)$ times. $\rightsquigarrow \cdots \rightsquigarrow \text{Time } \mathcal{O}\left(\frac{n}{p}(\log n + \beta \log p) + \alpha \log^2 p\right)$



Distributed memory parallel quicksort

Function parQuickSort(s : Sequence of Element, $i, j : \mathbb{N}$) : Sequence of Element

$$\begin{array}{l} p' \coloneqq j - i + 1 \\ \textbf{if } p' = 1 \textbf{ then } \texttt{quickSort}(s) \quad ; \textbf{return } \texttt{s} \qquad // \texttt{ sort locally} \\ v \coloneqq \texttt{pickPivot}(s, i, j) \\ a \coloneqq \langle e \in s : e \leq v \rangle; \quad b \coloneqq \langle e \in s : e > v \rangle \\ n_a \coloneqq \sum_{i \leq k \leq j} |a| @k; \quad n_b \coloneqq \sum_{i \leq k \leq j} |b| @k \\ k' \coloneqq \frac{n_a}{n_a + n_b} p' \\ \texttt{choose } k \in \{\lfloor k' \rfloor, \lceil k' \rceil\} \texttt{ such that } \max\left\{ \lceil \frac{n_a}{k} \rceil, \lceil \frac{n_b}{p' - k} \rceil \right\} \texttt{ is minimized} \\ \texttt{send the } a \texttt{-s to PEs } i..i + k - 1 (\leq \lceil \frac{n_a}{k} \rceil \texttt{ per PE}) \\ \texttt{send the } b \texttt{-s to PEs } i + k..j (\leq \lceil \frac{n_b}{p' - k} \rceil \texttt{ per PE}) \\ \texttt{receive data sent to PE } i_{PE} \texttt{ into } s \\ \textbf{if } i_{PE} < i + k \texttt{ then } \texttt{parQuickSort}(s, i, i + k - 1) \texttt{ else } \texttt{parQuickSort}(s, i + k, j) \end{array}$$



140



Load Balance

Simplified scenario: splitting always with ratio 1:2 larger subproblem gets one PE-load too much. Imbalance-factor:

$$\begin{split} \prod_{i=1}^{k} 1 + \frac{1}{p\left(\frac{2}{3}\right)^{i}} &= e^{\sum_{i=1}^{k} \ln\left(1 + \frac{1}{p\left(\frac{2}{3}\right)^{i}}\right)} \\ &\leq e^{\sum_{i=1}^{k} \frac{1}{p\left(\frac{2}{3}\right)^{i}}} = e^{\frac{1}{p}\sum_{i=0}^{k} \left(\frac{3}{2}\right)^{i}} \text{ geom. sum} \\ &= e^{\frac{1}{p} \frac{\left(\frac{3}{2}\right)^{k+1} - 1}{\frac{3}{2} - 1}} \leq e^{\frac{1}{p} 3\left(\frac{3}{2}\right)^{k}} = e^{3} \approx 20.1 \end{split}$$



Time
$$\mathscr{O}\left(\frac{n}{p}\log\frac{n}{p} + \log^2 p\right)$$



Better Balance?

Janus-quicksort? Axtmann, Wiebigke, Sanders, IPDPS 2018

 \Box for small p' choose pivot carefully

 \Box for small p' ($\Theta(\log p)$) switch to sample sort?

Alternative: always halve PEs, randomizaztion, careful choice of pivot Axtmann, Sanders, ALENEX 2017






Multi-Pivot Methods

Simplifying assumption: perfect splitters are available for free

```
//Für 0 < k < p let v_k the element with rank k \cdot n/p
// Set v_0 = -\infty and v_p = \infty.
initialize p empty messages N_k, (0 \le k < p)
for i := 0 to n/p - 1 do
    determine k, such that v_k < d_i \leq v_{k+1}
    put d_i into message N_k
                                                            // All-to-all
send N_i to PE i and
                                      II personalized communication
receive p messages
sort received data
```



Analysis



Idealizing assumption is realistic for permutation.

147

Sample Sort

choose a total of *ap* random elements s_k , (*a* per PE) ($1 \le k \le ap$) sort $[s_1, \ldots, s_{ap}]$ // or only for i := 1 to p - 1 do $v_i := s_{ai}$ // multiple selection

$$v_0:=-\infty; \quad v_P:=\infty$$



148



Lemma 2. $a = \mathcal{O}\left(\frac{\log n}{\varepsilon^2}\right)$ suffices such that with probability $\geq 1 - \frac{1}{n}$ no PE gets more than $(1 + \varepsilon)n/p$ elements.



Lemma: $a = \mathscr{O}\left(\frac{\log n}{\varepsilon^2}\right)$ suffices such that with probability $\geq 1 - \frac{1}{n}$ no PE gets more than $(1 + \varepsilon)n/p$ elements.

Proof idea: We analyze and algorithm that choses global samples with replacement.

Let $\langle e_1, \ldots, e_n \rangle$ denote the input in sorted order. fail: Some PE gets more than $(1 + \varepsilon)n/p$ elements $\rightarrow \exists j :\leq a \text{ samples from } \langle e_j, \dots, e_{j+(1+\varepsilon)n/p} \rangle \text{ (event } \mathscr{E}_j)$ $\rightarrow \mathbb{P}[\text{fail}] \leq n \mathbb{P}[\mathscr{E}_i], j \text{ fixed.}$ Let $X_i := \begin{cases} 1 & \text{if } s_i \in \langle e_j, \dots, e_{j+(1+\varepsilon)n/p} \rangle \\ 0 & \text{else} \end{cases}$, $X := \sum_i X_i$ $\mathbb{P}\left[\mathscr{E}_{j}\right] = \mathbb{P}\left[X < a\right] = \mathbb{P}\left[X < 1/(1+\varepsilon)\mathbb{E}[X]\right] \approx \mathbb{P}\left[X < (1-\varepsilon)\mathbb{E}[X]\right]$ $\mathbb{E}[X_i] = \mathbb{P}[X_i = 1] = \frac{1+\varepsilon}{n}$



٠

Chernoff-Bound

Lemma 3. Let $X = \sum_{i} X_{i}$ denote the sum of independent 0-1 random variables.

$$\mathbb{P}[X < (1 - \varepsilon)\mathbb{E}[X]] \le \exp\left(-\frac{\varepsilon^2\mathbb{E}[X]}{2}\right)$$

Applied to our problem:

$$\mathbb{P}[X < a] \le \exp\left(-\frac{\varepsilon^2(1+\varepsilon)a}{2}\right) \le \exp\left(-\frac{\varepsilon^2 a}{2}\right) \le \frac{1}{n^2}$$
$$\Leftrightarrow a \ge \frac{4}{\varepsilon^2} \ln n$$



Analysis of Sample Sort





Sorting Samples

- Using gather/gossiping
- Using gather-merge
- ☐ Fast ranking
- Parallel quicksort
- Recursively using sample sort

Sorting Samples

Using gather/gossiping

Using gather-merge

Fast ranking

Parallel quicksort

Recursively using sample sort





MPI Sample Sort – Init and Local Sample

Many thanks to Michael Axtmann

template <class element=""></class>	1
void parallelSort(MPI_Comm comm, vector <element>& data,</element>	2
MPI_Datatype mpiType, int p, int myRank)	3
{ random_device rd;	4
mt19937 rndEngine(rd());	5
uniform_int_distribution <size_t> dataGen(0, data.size() $- 1$);</size_t>	6
vector <element> locS; // local sample of elements from input <data></data></element>	7
const int a = (int)(16*log(p)/log(2.)); // oversampling ratio	8
for (size_t i=0; i < (size_t)(a+1); ++i)	9
locS.push_back(data[dataGen(rndEngine)]);	10

Find Splitters



vector <element> s(locS.size() * p); // global samples</element>	1
MPI_Allgather(locS.data(), locS.size(), mpiType,	2
s.data(), locS.size(), mpiType, comm);	3

```
sort(s.begin(), s.end()); // sort global sample 5
for (size_t i=0; i < p-1; ++i) s[i] = s[(a+1) * (i+1)]; //select splitters 6
s.resize(p-1); 7
```

Partition Locally



```
vector<vector<Element>> buckets(p); // partition data 1
for(auto& bucket : buckets) bucket.reserve((data.size() / p) * 2); 2
for( auto& el : data) { 3
    const auto bound = upper_bound(s.begin(), s.end(), el); 4
    buckets[bound - s.begin()].push_back(el); 5
}
data.clear(); 7
```

Find Message Sizes



```
// exchange bucket sizes and calculate send/recv information
                                                                            1
vector<int> sCounts, sDispls, rCounts(p), rDispls(p + 1);
                                                                            2
sDispls.push_back(0);
                                                                            3
for (auto& bucket : buckets) {
                                                                            4
  data.insert(data.end(), bucket.begin(), bucket.end());
                                                                            5
  sCounts.push_back(bucket.size());
                                                                            6
  sDispls.push_back(bucket.size() + sDispls.back());
                                                                            7
}
                                                                            8
MPI_Alltoall(sCounts.data(),1,MPI_INT,rCounts.data(),1,MPI_INT,comm); 9
// exclusive prefix sum of recv displacements
                                                                            10
rDispls[0] = 0;
                                                                            11
for(int i = 1; i \le p; i++) rDispls[i] = rCounts[i-1]+rDispls[i-1];
                                                                            12
```

Data Exchange and Local Sorting



```
vector<Element> rData(rDispls.back()); // data exchange 1
MPI_Alltoallv(data.data(), sCounts.data(), sDispls.data(), mpiType, 2
rData.data(), rCounts.data(), rDispls.data(), mpiType, comm); 3
```

```
sort(rData.begin(), rData.end());
rData.swap(data);
}
```



Experiments Speedup on $4 \times$ Intel E7-8890 v3



input size (elements per thread)



Sorting by multiway merging

Function mmSort(d, n, p) // shared memory not SPMD PE *i* sorts d[in/p..(i+1)n/p]; barrier synchronization PE *i* finds v_i with rank in/p in *d*; barrier synchronization PE *i* merges *p* subsequences with $v_k \le d_j < v_{k+1}$





Multisequence Selection

Idea: each PE determines a splitter with appropriate global Rank (shared memory)

Comparison based lower bound: $\Omega\left(p\log\frac{n}{p}\right)$

We present an algorithmus with $\mathscr{O}\left(p\log n\log \frac{n}{p}\right)$



Splitter Selection

Processor *i* selects the element with global rank $k = \frac{in}{p}$.

Simple algorithm: quickSelect exploiting sortedness of the sequences.



Idea:

Ordinary select but $p \times$ binary search instead of partitioning

Function msSelect(S : Array of Sequence of Element; $k : \mathbb{N}$) : Array of \mathbb{N}

for i := 1 to |S| do $(\ell_i, r_i) := (0, |S_i|)$ invariant $\forall i : \ell_i ... r_i$ contains the splitting position of S_i invariant $\forall i, j : \forall a \le \ell_i, b > r_j : S_i[a] \le S_j[b]$ while $\exists i : \ell_i < r_i$ do $v := \text{pickPivot}(S, \ell, r)$ for i := 1 to |S| do $m_i := \text{binarySearch}(v, S_i[\ell_i ... r_i])$ if $\sum_i m_i \ge k$ then r := m else $\ell := m$ return ℓ





Analysis of *p***-way mergesort**

$$T_{\text{pMergeSort}}(p,n) = \mathscr{O}\left(\underbrace{\frac{n}{p}\log\frac{n}{p}}_{\text{local sort}} + \frac{p\log n\log\frac{n}{p}}{\frac{p}{\text{ms-selection}}} + \underbrace{\frac{n}{p}\log p}_{\text{merging}}\right)$$

 \Box efficient if $n \gg p^2 \log p$

deterministic (almost)

perfect load balance

somewhat worse constant factors than sample sort



Distributed Multisequence Selection

Owner computes paradigm

 $\mathcal{O}(\log n)$ global levels of recursion.

Gather + Broadcast for finding pivot / distribution (vector length p-1).

p-1 local searches everywhere.

Reduction for finding partition sizes (Vector length p-1).

Expected time
$$\mathscr{O}\left(\log n\left(p(\log \frac{n}{p} + \beta) + \log p\alpha\right)\right)$$



Distributed Multisequence Selection

Function dmSelect(s : Seq of Elem; k : Array[1..p] of \mathbb{N}) : Array[1..p] of \mathbb{N}

 ℓ, r, m, v, σ : Array [1..p] of N for i := 1 to p do $(\ell_i, r_i) := (0, |s|)$ // initial search ranges while $\exists i, j : \ell_i @ j \neq r_i @ j$ do // or-reduction $v := pickPivotVector(s, \ell, r)$ // reduction, prefix sum, broadcast for i := 1 to p do $m_i := binarySearch(v_i, s[\ell_i..r_i])$ $\sigma := \sum_i m@i$ // vector valued reduction for i := 1 to p do if $\sigma_i \ge k_i$ then $r_i := m_i$ else $\ell_i := m_i$ return ℓ



CRCW Sorting in logarithmic time

Consider case n = p.

- \Box sample of size \sqrt{p}
- $\Box k = \Theta(\sqrt{p}/\log p)$ splitters
 - \Box Buckets have size $\leq cp/k$ elements whp
 - Allocate buckets of size 2cp/k
 - Write elements to random free positions within their bucket
 - Compactify using prefix sums

Recursion



Example





More on Sorting I

Cole's merge sort: [JáJá Section 4.3.2] Time $\mathscr{O}\left(\frac{n}{p} + \log p\right)$ deterministic, EREW PRAM (CREW in [JáJá]). Idea: Pipelined parallel merge sort. Use (deterministic) sampling to predict where data comes from.

Sorting Networks: nodes sort 2 elements. Simple networks have $\mathscr{O}(\log^2 n)$ depth (e.g. bitonic sort). They yield reasonable deterministic sorting algorithms (2 elements \rightsquigarrow merge-and-split of two subsequences). Very complicated ones with depth $\mathscr{O}(\log n)$.



More on Sorting II

Integer Sorting: (Close to) linear work. Very fast algorithms on CRCW PRAM.

Multi-Pass-Sample/Merge-Sort: more general compromise between latency and communication volume, e.g. AMS-Sort Axtmann, Bingmann, Schulz, Sanders SPAA 2015



172



Slowdown wrt Fastest Algorithm





Collective Communication

Broadcast

Reduction

Prefix sum

- Not here hier: gather / scatter
- Gossiping (= all-gather = gather + broadcast)
- All-to-all personalized communication
 - equal message lengths
 - arbitrary message lengths, = h-relation



Prefix sums

[Leighton 1.2.2] Compute

$$x@i:=\bigotimes_{i'\leq i}m@i'$$

(on PE i, m may be a vector with n bytes length.)





Plain Pipeline

As in broadcast





001

a-b

a–b

011

c-d

c–d

101

e-f

e-f

111

g-h

Hypercube Algorithm

//view PE index *i* as a *lld*-bit bit array **Function** hcPrefix(m)

> $x := \sigma := m$ for k := 0 to d - 1 do invariant $\sigma = \bigotimes_{\substack{i = i [k..d-1]0^k}}^{i[k..d-1]1^{\kappa}} m@j$ **invariant** $x = \bigotimes_{j=i[k..d-1]0^k}^{i} m@j$ $y := \sigma @(i \oplus 2^k) \qquad \qquad // \text{ sendRecv}$ $\sigma := \sigma \otimes y$ if i[k] = 1 then $x := y \otimes x$

100

e-e

e-e

110

g-g

g-g

000

a–a

a–a

010

с-с

с-с

101

f-f

f-f

111

h–h

Х

001

b-b

b-b

011

d–d

d-d

100

e-f

e-e

110

g-h

000

a-b

a-a

010

c-d

с-с





return x



Analysis

Telephone model:

$$T_{\text{prefix}} = (\alpha + n\beta) \log p$$

Pipelining does not work since all PEs are busy.

Pipelined Binary Tree Prefix Sums

Inorder numbering of the nodes

Upward phase: as with reduction but PE *i* stores $\sum_{j=i'}^{i} x@j$

Downward phase: PE *i* receives $\sum_{j=1}^{i'-1} x@j$

(root: = 0 !)

and forwards this to the left. right subtree gets $\sum_{i=1}^{i} x@j$

Each PE only active once per phase. \rightarrow pipelining OK



Pseudocode

Function InOrderTree::prefixSum(m)

// upward phase: x := 0; receive(leftChild, x) z := 0; receive(rightChild, z) send(parent, x + m + z)

//downward phase: $\ell := 0;$ receive(parent, ℓ) send(leftChild, ℓ) send(rightChild, $\ell + x + m$)

return $\ell + x + m$






23-Prefix Sums

Numbering is inorder-numbering for both trees!



Analysis







Generalization:

Applies to any algorithm based on inorder numbered trees

 \rightsquigarrow ESBT does not work?



Each PE has a message m of length n.

At the end, each PE should know all messages.

Hypercube Algorithm

Let ' \cdot ' denote the concatenation operation; $p=2^d$



$$y := m$$

for $0 \le j < d$ do
 $y' := \text{the } y \text{ from PE } i \oplus 2^j$
 $y := y \cdot y'$
return y





Example





Analysis

Telephone model, $p = 2^d$ PEs, *n* byte per PE:

$$T_{\text{gossip}}(n,p) \approx \sum_{j=0}^{d-1} \alpha + n \cdot 2^j \beta = \log p \alpha + (p-1)n \beta$$

All-Reduce

Reduction instead of concatenation.

Advantage: Factor two less startups than reduction plus broadcast

Disadvantage: $p \log p$ messages.

This a disadvantage for congenstion-prone networks.



All-to-all Personalized Communication

Each PE has p - 1 messages of length n. One for each other PE. m[i]@i is for PE i itself

Hypercube Algorithm

 $\mathsf{PE}\ i$

```
for j := d - 1 downto 0 do

Get from PE i \oplus 2^j all its messages

destined for my j-D subcube

Move to PE i \oplus 2^j all my messages

destined for its j-D subcube
```



Analysis, Telephone model:

$$T_{\text{all-to-all}}(p,n) \approx \log p(\frac{p}{2}n\beta + \alpha)$$

Fully Connected:

When n is large, rather send messages individually (Factor $\log p$ less communication volume)

1-Factor-Algorithm

[König 1936]

p odd, i is PE-index:

for
$$r := 0$$
 to $p - 1$ do
 $k := 2r \mod p$
 $j := (k - i) \mod p'$
 $send(j, m_{ij}) || recv(j, m_{ji})$

pairwise communication (telephone model): The partner of the partners of j in round i is $i - (i - j) \equiv j \mod p$ Time: $p(n\beta + \alpha)$ optimal for $n \to \infty$





Sanders: Parallel Algorithms January 17, 2022

1-Factor Algorithm

p even:

// PE index $j \in \{0, ..., p-1\}$ for i := 0 to p - 2 do idle:= $\frac{p}{2}i \mod (p-1)$ if j = p - 1 then exchange data with PE idle else

if j = idle then exchange data with PE p - 1else exchange data with PE $(i - j) \mod (p - 1)$

Time: $(p-1)(n\beta + \alpha)$ optimal für $n \to \infty$







Data exchange with irregular message lengths

- \Box Particularly interesting with all-to-all \rightarrow sorting
- similar problems in inhomogeneous interconnection networks or competition by other jobs.

Sanders: Parallel Algorithms January 17, 2022

The "Ostrich"-Algorithm

Push all messages using asynchronous send operations.

Receive what comes in.

Ostrich-Analysis:

BSP-Model: Time L + gh

But what is L and g in our single-ported models?



h-**Relation**

 $h_{in}(i) := \#$ packets received by PE *i* $h_{\text{out}}(i) := \#$ packets sent by PE *i*

simplex:
$$h := \max_{i=1}^{p} h_{in}(i) + h_{out}(i)$$

duplex: $h := \max_{i=1}^{p} \max(h_{in}(i), h_{out}(i))$

Lower bound for packet-wise delivery: h steps, i.e., Time $h(\alpha + |\text{packet}|\beta)$



Offline *h***-Relations in the Duplex Model**₁

[König 1916]

Consider the bipartite multigraph

 $G = (\{s_1, \dots, s_p\} \cup \{r_1, \dots, r_p\}, E) \text{ with} \\ |\{(s_i, r_j) \in E\}| = \# \text{ packets from PE } i \text{ to PE } j.$

Theorem: \exists edge coloring $\phi : E \rightarrow \{1..h\}$, i.e., no two equal-colored edges are incident to any node.

for *j* := 1 **to** *h* **do**

send messages with color j

optimal when postulating packet-wise delivery







Offline *h***-Relations in the Duplex Model**

Problems:

- Computing edge colorings online
 is complicated and expensive
 - Shredding messages into packets increases # startups







Offline *h***-Relations in the Simplex-Model**

[Petersen 1891? Shannon 1949?]

Consider the multigraph $G = (\{1, ..., p\}, E)$ with $|\{\{i, j\} \in E\}| = \#$ packets between PE *i* and PE *j* (both directions).

Theorem: \exists edge coloring $\phi : E \rightarrow \{1..3 \lfloor h/2 \rfloor + h \mod 2\}$

for j := 1 to h do Send messages of color j

optimal???





٠

How Helper Hasten *h***-Relations**

[Sanders Solis-Oba 2000]

Satz 4. For h-relations in the simplex model,

$$\# \text{steps} = \begin{cases} \frac{6}{5}(h+1) & \text{if } p \text{ even} \\ (\frac{6}{5} + \frac{2}{p})(h+1) & \text{if } p \text{ odd} \end{cases}$$

On the other hand, there is a lower bound

$$\# \text{steps} \geq \begin{cases} \frac{6}{5}h & \text{if } p \text{ even} \\ (\frac{6}{5} + \frac{18}{25p})h & \text{if } p \text{ odd} \end{cases}$$



A Very Simple Case



Sanders: Parallel Algorithms January 17, 2022





Reduction *h***-Relation** $\rightsquigarrow \left| \frac{h}{2} \right|$ **2-Relations**

Karlsruhe Institute of Technology

200

Ignore direction of communications for now

- Connect nodes with odd degree
 - \rightsquigarrow all nodes have even degree
- Eulertour-technique: decompose the graph into edge disjoint cycles
 - Direct cycles in clockwise direction

 \rightsquigarrow indegree and outdegree $\leq \lceil h/2 \rceil$

- Build bipartite graph (as before)
- Color bipartite graph
- Color classes in bipartite graph → edge disjoint simple cycles in the input graph (2-relations)
- Reinstate orginal communication direction

Routing 2-Relations for Even *p*

201

Pair odd cycles.

1-cycles have nothing to do \rightsquigarrow most simple case



Two odd cycles with \geq 3 nodes

Split packets into 5 subpackets



Then turn this around



Odd p

Idea: Delete one edge in each 2-factor

Do that "always somewhere else"

Collect $\Theta(p)$ removed edges into a matching

 \rightsquigarrow one additional step every $\Theta(p)$ 2-factors.



Open Problems

- Get rid of splitting into 5 subpackages?
- Conjecture:

One *h*-Relation with $\leq \frac{3}{8}hp$ packets can be delivered in $\approx h$ steps.

- Explicitly account for startup overheads
- Explicitly account for interconnection network?
- Distributed scheduling



A Simple Distributed Algorithm — The Two-Phase-Algorithm

Idea: Irreg. All-to-all \rightarrow 2× regular All-to-all

Simplifying assumptions:

All message lengths are divisible by p (in doubt round up)

communication "with oneself" is accounted for

All PEs send and receive exactly *h* bytes
 (In doubt "pad" the messages)

Sanders: Parallel Algorithms January 17, 2022



||n[i]| is length of message m|i|**Procedure** alltoall2phase(m[1..p], n[1..p], p)for i := 1 to p do $a[i] := \langle \rangle$ for j := 1 to p do $a[i] := a[i] \odot m[j][(i-1)\frac{n[j]}{n} + 1..i\frac{n[j]}{n}]$ b := regularAIIToAII(a, h, p) $\delta := \langle 1, \ldots, 1 \rangle$ for i := 1 to p do $c[i] := \langle \rangle$ **for** *j* := 1 **to** *p* **do** $c[i] := c[i] \odot b[j] [\delta[j] .. \delta[j] + \frac{n[i]@j}{p} - 1]$ // Use All- $\delta[j] := \delta[j] + \frac{n[i]@j}{n}$ // gather to implement '@' d := regularAllToAll(c, h, p)permute d to obtain the desired output format





207



More on the Two-Phase-Algorithm

- □ Large p, small messages \rightsquigarrow split local data into $\mathscr{O}(p\log p)$ pieces (not p^2) and disperse randomly.
- Split the problems into regular and irregular part ~> two-phase protocol only applied to a part of the data.
 ~> open problem: how to split?



209

A Non-Preemptive Offline Algorithm (simplex)

[Sanders Solis-Oba 99, unpublished]

Goal: deliver all messages directly, as a whole.

Let $k := \max$. # messages one PE is involved in.

Time for executing the schedule $k\alpha + 2h\beta$. Here *h* is measured in bytes!



Abstract Description

- *s*:= empty schedule
- M:= set of messages to be scheduled

while $M \neq \emptyset$ do

 $t := \min \{ t : \exists m \in M : m \text{'s src and dest are idle at time } t \}$

 $s := s \cup$ "start sending *m* at time *t*"

 $M := M \setminus \{m\}$

Can be implemented such that per message, the time for $\mathcal{O}(1)$ priority-queue operationen and one *p*-bit bit-vector operation is used.

 \rightsquigarrow practicable for message lengths $\gg p$ and moderate p.



Open Problems for Non-Preemptive Offline Algorithms

- implement, measure, use, e.g. sorting, construction of suffix-arrays
- Better approximation algorithms?
- Parallel scheduling algorithms

Sanders: Parallel Algorithms January 17, 2022



Summary: All-to-All

Ostrich: Delegate to online, asynchronous routing.

Good when that is implemented well.

Regular+2Phase: more robust. But, factor 2 is troubling. A lot of copying overhead.

Non-preemptive: Minimizes startups and communication volume. Faktor2 (worst case). Centralized Scheduling is troubling.Good for repeated identical problems.

Coloring-based algorithms: Almost optimal for large packets. Complex. Distributed Implementation? Splitting into packets is troubling.

Comparison of approaches?

Parallel Priority Queues



Manage a set M of elements. n = |M|. Initially empty

Binary Heaps (sequential)

Procedure insert(e) $M := M \cup \{e\}$ // $\mathcal{O}(\log n)$ Function deleteMin $e := \min M;$ $M := M \setminus \{e\};$ return $e// \mathcal{O}(\log n)$



Parallel Priority Queues, Goal

insert*: Each PE inserts a constant number of elements, time $\mathcal{O}(\log n + \log p)$?

deleteMin*: delete the p smallest elements, time $\mathcal{O}(\log n + \log p)$?

Asynchronous Variant (later): Each PE can insert or delete at any time. Semantics: \exists temporal ordering of the operations consistent with sequential execution.



Applications

- Priority-driven scheduling
- Best first branch-and-bound:

Find best solution in a large, implicitly defined tree. (later more)

Discrete event simulation

Naive Implementation

PE 0 manages a sequential PQ All others send requests

insert: $\Omega(p(\alpha + \log n))$

deleteMin: $\Omega(p(\alpha + \log n))$




Branch-and-Bound

H: tree (V, E) with bounded node degree

c(v): node costs — growing when descending a path

- v^* : leaf with minimal costs
- $\tilde{V}: \{ v \in V : v \le v^* \}$
- *m*: $|\tilde{V}|$ Simplification: $\Omega(p \log p)$
- *h*: Depth of \tilde{H} (subgraph of *H* induced by \tilde{V}).
- $T_{\rm x}$: Time for generating the successors of a node
- T_{coll} : Upper bound for broadcast, min-reduction, prefix-sum, routing one element from/to random partner. $\mathscr{O}(\alpha \log n)$ on many networks
 - $\mathscr{O}(\alpha \log p)$ on many networks.

Sequential Branch-and-Bound



 $\begin{array}{ll} Q = \{ {\rm root \ node} \} : {\rm PriorityQueue} & // \ {\rm frontier \ set} \\ c^* = \infty & // \ {\rm best \ solution \ so \ far} \\ {\rm \textbf{while}} \ Q \neq \emptyset \ {\rm \textbf{do}} \\ v := Q. {\rm deleteMin} \\ {\rm \textbf{if}} \ c(v) < c^* \ {\rm \textbf{then}} \\ {\rm \textbf{if}} \ v \ {\rm is \ a \ leaf \ node \ \textbf{then}} \ {\rm process \ new \ solution; \ } c^* := c(v) \\ {\rm \textbf{else} \ insert \ successors \ of \ } v \ {\rm into} \ Q \end{array}$

 $T_{\text{seq}} = m(T_{\text{x}} + \mathcal{O}(\log m))$

Parallel Branch-and-Bound

 $Q = \{\text{root node}\}$: **Paralle**|**PriorityQueue**] $c^* = \infty$ while $Q \neq \emptyset$ do v := Q.deleteMin*if $c(v) < c^*$ then if v is a leaf node then process new solution update c^* else insert successors of v into Q



```
// best solution so far
```

// SPMD!

// Reduction

Analysis

Theorem:
$$T_{\text{par}} = \left(\frac{m}{p} + h\right) \left(T_{\text{x}} + \mathscr{O}\left(T_{\text{queueOp}}\right)\right)$$

Case 1:

(at most m/p lterations): All processed nodes are in \tilde{V}

Case 2:

(at most h Iterations): Some nodes outside \tilde{V} are processed \rightarrow maximal path length from a node in Qto the optimal solutions is being reduced.





The Karp–Zhang Algorithm

 $Q = \{\text{root node}\}$: PriorityQueue // local! $c^* = \infty$ // best solution so far while $\exists i : Q@i \neq \emptyset$ do // local! v := Q.deleteMin if $c(v) < c^*$ then if v is a leaf node then process new solution $c^* := \min_i c(v) @i$ // Reduction else for each successor v' of v do insert v into Q@i for random i

Satz: Expected time is asymptotically optimal





Our Approach PE: 2 3 4 1 **B&B** Processes New Nodes Random Placement Local Queues Top-Nodes A Filter p best Assign to PEs



Parallel Priority Queues: Approach

The queue is the union of local queues

Insert sends new elements to random local queues Intuition: each PE gets a representative view on the data

deleteMin* finds the globally smallest elements

(act locally think globally) PE: 1 2 3 4 distributes them to the PEs Filter p best Assign to PEs



Simple Probabilistic Properties

With high probability (whp):

here $\geq 1 - p^{-c}$ for a constant c we can choose

 \Box whp only $\mathscr{O}\left(\frac{\log p}{\log \log p}\right)$ elements per local queue during insertion

 \Box whp, the $\mathscr{O}(\log p)$ smallest elements of each local queue together contain the *p* globally best elements

 \Box whp no local queue contains more then $\mathscr{O}(n/p + \log p)$ elements

Proof: Chernoff-bounds again-and-again.

(standard setting, balls-into-bins)



Parallel Implementation I

Insert

Sending: T_{coll} Local insertions: $\mathscr{O}\left(\frac{\log p}{\log \log p} \cdot \log \frac{n}{p}\right)$. (Better with "advanced" local queues Careful: amortized bounds are not sufficient.)

Parallel Implementation I

deleteMin*

Procedure deleteMin^{*}(Q_1, p) Q_0 := the $\mathcal{O}(\log p)$ smallest elements of Q_1 M:= select(Q_0, p) // later enumerate $M = \{e_1, \dots, e_p\}$ assign e_i to PE i // use prefix sums **if** max_i $e_i > \min_j Q_1 @ j$ **then** expensive special case treatment empty Q_0 back into Q_1



Analysis

Remove locally: $\mathscr{O}\left(\log p \log \frac{n}{p}\right)$ Selection: $\mathscr{O}(T_{\text{coll}})$ whp todo Enumerate M: $\mathscr{O}(T_{\text{coll}})$ Deliver results: $\mathscr{O}(T_{\text{coll}})$ (random sources) Verify: $\mathscr{O}(T_{\text{coll}}) +$

(sth polynomial in p)·(a polynomially small probability)

Insert locally: $\mathscr{O}\left(\log p \log \frac{n}{p}\right)$





Parallel Implementation II

Idea: Avoid ping-pong of the $\mathcal{O}(\log n)$ smallest elements.

Split the queue into Q_0 and Q_1 . Invariant: whp $|Q_0| = \mathscr{O}(\log p)$





Parallel Implementation II

Insert

Send: T_{coll}

Insert locally: merge Q_0 and new elements $\mathscr{O}(\log p)$ whp.

```
Cleanup: Empty Q_0 every \log p iterations.

\rightsquigarrow \cos \mathscr{O}\left(\log p \log \frac{n}{p}\right) per \log p iterations

\rightsquigarrow \text{average costs } \mathscr{O}\left(\log \frac{n}{p}\right)
```

Parallel Implementation II

deleteMin*

Procedure deleteMin* (Q_0, Q_1, p) while $|\{e \in \breve{Q}_0 : e < \min\breve{Q}_1\}| < p$ do $Q_0 := Q_0 \cup \{\text{deleteMin}(Q_1)\}$ $M := \text{select}(Q_0, p)$ // later enumerate $M = \{e_1, \dots, e_p\}$ assign e_i to PE i // use prefix sums



Analysis

Remove locally:
$$\mathscr{O}(1)$$
 expected iterations $\rightsquigarrow \mathscr{O}\left(T_{\text{coll}} + \log \frac{n}{p}\right)$

Selection: $\mathcal{O}(T_{\text{coll}})$ whp todo

Enumerate *M*: $\mathcal{O}(T_{\text{coll}})$

Deliver results: $\mathcal{O}(T_{\text{coll}})$ (random sources)



Sanders: Parallel Algorithms January 17, 2022

Result

insert*: expected
$$\mathscr{O}\left(T_{\text{coll}} + \log \frac{n}{p}\right)$$

deleteMin*: expected $\mathscr{O}\left(T_{\text{coll}} + \log \frac{n}{p}\right)$





Randomized Selection [Blum et al. 1972]

Given n (randomly allocated) elements Q, find the k smallest ones.

 \Box choose a sample ${f s}$

$$\exists u := \text{ element with rank } \frac{k}{n} |\mathbf{s}| + \Delta \text{ in } \mathbf{s}.$$

 $\ell := \text{ element with rank } \frac{k}{n} |\mathbf{s}| - \Delta \text{ in } \mathbf{s}.$

] Partition Q into

 $Q_{<} := \{q \in Q : q < \ell\},\$ $Q_{>} := \{q \in Q : q > u\},\$ $Q' := Q \setminus Q_{<} \setminus Q_{>}$

 $\Box \ \ {\rm If} \ |Q_<| < k \ {\rm and} \ |Q_<| + |Q'| \ge k, \ {\rm output} \ Q_< \ {\rm and} \ {\rm find \ the} \\ k - |Q_<| \ {\rm smallest \ elements \ of} \ Q'$

] All other cases are unlikely if $|\mathbf{s}|, \Delta$ are sufficiently large.



Randomized Selection [Blum et al. 1972]





Parallel Implementation

$$\square |\mathbf{s}| = \sqrt{p} \rightsquigarrow$$
 sample can be sorted in time $\mathscr{O}(T_{\text{coll}})$.

$$\Box \ \Delta = \Theta\left(p^{1/4+\varepsilon}\right) \text{ for a small constant } \varepsilon.$$

This makes difficult cases unlikely.

No elements are redistributed. Random initial distribution guarantees good load balance whp.

 \Box Whp a constant number of iterations suffices until only \sqrt{p} elements are left. \leadsto Then sort directly.

Overall expected time $\mathscr{O}\left(\frac{n}{p} + T_{\text{coll}}\right)$



Parallel Priority Queues – Refinements

Procedure deleteMin* (Q_0, Q_1, p) while $|\{e \in \check{Q}_0 : e < \min\check{Q}_1\}| < p$ do $Q_0 := Q_0 \cup \{\text{deleteMin}(Q_1)\}$ // select immediately $M := \text{select}(Q_0, p)$ // laterenumerate $M = \{e_1, \dots, e_p\}$ // use prefix sums

Or just use sufficiently many locally smallest elements and check later



Parallel Priority Queues – Refinements

Mergable priority queues?

Bulk delete after flush?

Larger samples

Remove larger batches?

Only a subset of the PEs work as PQ-server?

Selection by pruned merging:

A reduction with vector length $\mathscr{O}(\sqrt{p\log p})$

Sanders: Parallel Algorithms January 17, 2022

Asynchronous Variant

238

Accept insertions but do not immediately insert.

Batched deleteMin in a buffer.

Access buffer using an asynchronous FIFO.

Sometimes:

Invalidate FIFO,

commit inserted elements

refill buffer



Implementation on IBM SP-2, $m = 2^{24}$





Implementation on Cray T3D, $m = 2^{24}$

p = 256

insert 256 elements and a deleteMin*: centralized: > 28.16ms parallel: 3.73ms

break-even at 34 PEs

Sanders: Parallel Algorithms January 17, 2022



More on parallel Priority Queues – History

Different approach starts with a binary heap:

nodes with p sorted Elements.

Invariant: All elements > all elements in parent node

Compare-and-swap \rightsquigarrow merge-and-split

Elegant but expensive

Parallelization of a single access \rightsquigarrow constant time with $\log n$ PEs.



Communication Efficient Priority Queues

Each PE stores a search tree augmented with subtree sizes. ~>>

- \Box local insert $\mathscr{O}(\log n)$ time.
-] find k smallest elements in time $\mathscr{O}(\log^2 n)$

(similar to multi-sequence selection for mergesort)

 \Box find $\Theta(k)$ smallest elements in time $\mathscr{O}(\log n)$

Communication Efficient Algorithms for Top-k Selection Problems, with Lorenz Hübschle-Schneider, IPDPS 2016



MultiQueues: Simple Relaxed Concurrent

With Roman Dementiev and Hamza Rihani, Marvin Williams SPAA 2015, ESA 2021

 $\Box cp$ local queues $Q[1], \ldots, Q[cp]$, constant c > 1, e.g., c = 2

Insert into random local queues ("wait-free" locking)

Delete smallest elements from c randomly chosen queues

Improve locality using insertion/deletion buffers

Optional: Stickyness – stick to the same set of queues for s consecutive operations

244

Throughput





Rank Error

Average rank of deleted elements



Delay



Average number of larger elements deleted before a removed element





List Ranking

Motivation:

With arrays a[1..n] we can do many things in parallel

 \Box PE *i* processes $a[(i-1)\frac{n}{p}+1..i\frac{n}{p}]$

Prefix sums

. . . .

Can we do the same with linked lists?

Yes! For example, convert into an array

Sanders: Parallel Algorithms January 17, 2022

List Ranking

L: List

n: Elements

S(i): Successor of element i(unordered) S(i) = i: at end of list



248

P(i): Predecessor of element i

Exercise: compute in constant time for n PE PRAM

R(i): Rank. Initially 1, 0 for last element.

Output: R(i) = distance to the end (when following the chain).

Array-Conversion: store item *i* in a(n - R(i))



Motivation II

Lists are very simple graphs

- \rightsquigarrow warmup for graph algorithms
- \rightsquigarrow long paths are a main obstacle for parallelization.
- I.e., solutions might help for more general graphs also (at least trees)?

Pointer Chasing

find *i* such that S(i) = ifor r := 0 to n - 1 do R(i) := r i := P(i) \Box Work $\mathcal{O}(n)$ \Box Time $\Theta(n)$



// parallelizable

// inherently sequential?



Doubling using CREW PRAM, n = p

$$\begin{split} Q(i) &\coloneqq S(i) \qquad // \text{ SPMD. PE index } i \\ \text{invariant } \sum_{j \in Q_i} R(j) = \text{rank of item } i \\ Q_i \text{ is the positions given by} \\ \text{chasing } Q \text{-pointers from pos } i \\ \text{while } R(Q(i)) \neq 0 \text{ do} \\ R(i) &\coloneqq R(i) + R(Q(i)) \\ Q(i) &\coloneqq Q(Q(i)) \end{split}$$

Analysis

Induction Hypothesis: After k iterations

- $\square R(i) = 2^k$ or
- $\square R(i) =$ final result
- Proof: True for k = 0. $k \rightsquigarrow k+1$:

Case $R(i) < 2^k$: already final value Case $R(i) = 2^k$, $R(Q(i)) < 2^k$: now final value (Invariant, IH) Case $R(i) = R(Q(i)) = 2^k$: Now 2^{k+1}

Work $\Theta(n \log n)$

Time $\Theta(\log n)$




Independent Set Removal

//Compute the sum of the R(i)-values when following the S(i) pointers **Procedure** independentSetRemovalRank(n, S, P, R)

if $p \ge n$ then use doubling; return find $I \subseteq 1..n$ such that $\forall i \in I : S(i) \notin I \land P(i) \notin I$ find a bijective mapping $f : \{1..n\} \setminus I \rightarrow 1..n - |I|$ foreach $i \notin I$ dopar // remove independent set I S'(f(i)) :=if $S(i) \in I$ then f(S(S(i))) else f(S(i)) $P'(f(i)) := \text{ if } P(i) \in I \text{ then } f(P(P(i))) \text{ else } f(P(i))$ R'(f(i)) :=**if** $S(i) \in I$ **then** R(i) + R(S(i)) **else** R(i)independentSetRemovalRank(n - |I|, S', P', R')foreach $i \notin I$ dopar R(i) := R'(f(i))foreach $i \in I$ dopar R(i) := R(i) + R'(f(S(i)))







Finding Independent Sets



Monte Carlo algorithm \rightsquigarrow Las Vegas algorithm: repeat until $|I| > \frac{n}{5}$. Expected time: $\mathcal{O}(n/p)$

Neither start nor end of list are in *I*.

Finding a Bijective Mapping

Prefix sum over the indicator function of $\{1..n\} \setminus I$:

$$f(i) = \sum_{j \le i} [j \notin I]$$







More on List Ranking



- Simple algorithm with expected time $\mathcal{O}(\log n)$
- \Box Complicated algorithm with worst case Time $\mathscr{O}(\log n)$
 - Many "applications" in PRAM algorithms
- Implementation on distributed-memory parallel computers [Sibeyn 97]: p = 100, $n = 10^8$, speedup 30.
- Generalization for segmented lists, trees
- Generalization for general graphs:
 contract nodes or edges
- Example for multilevel algorithms



Newer Implementation Results

\Box Cut the list at s random place

Sequential algorithm for each sublist

 \Box Recursive solution on instance of size *s*

Speedup ≈ 10 on 8-core CPU (???) [Wei JaJa 2010]



Parallel Graph Algorithms

The "canonical" "easy" graph algorithms:

Main interest, sparse, polylog. execution time, efficient

- DFS
- BFS
- Shortest paths

(nonnegative SSSP $\mathscr{O}(n)$ par. time. interesting for $m = \Omega(np)$) (what about APSP?)

- Topological sorting
- + Connected components (but not strongly connected)
- + Minimal spanning trees
- + Graph partitioning

Minimum Spanning Trees

Undirected graph G = (V, E)Nodes V, n = |V|, e.g., $V = \{1, ..., n\}$ Edges $e \in E, m = |E|$, two-element subsets of VEdge weight $c(e), c(e) \in \mathbb{R}_+$ wlog all different G is connected, i.e., \exists path between any two nodes.



Find a tree (V, T) with minimum weight $\sum_{e \in T} c(e)$ that connects all nodes.





Selecting and Discarding MST Edges

The Cut Property

For any $S \subset V$ consider the cut edges $C = \{\{u, v\} \in E : u \in S, v \in V \setminus S\}$

The lightest edge in C can be used in an MST.



The Cycle Property

The heaviest edge on a cycle is not needed for an MST



The Jarník-Prim Algorithm

[Jarník 1930, Prim 1957]

Idea: grow a tree

 $T := \emptyset$ $S := \{s\} \text{ for arbitrary start node } s$ $repeat \ n - 1 \text{ times}$ find (u, v) fulfilling the cut property for S $S := S \cup \{v\}$ $T := T \cup \{(u, v)\}$







Graph Representation for Jarník-Prim Adjacency Array

We need node \rightarrow incident edges



Analysis

 $\square \ \mathscr{O}(m+n)$ time outside priority queue

 \square *n* deleteMin (time $\mathcal{O}(n \log n)$)

 $\square \ \mathscr{O}(m)$ decreaseKey (time $\mathscr{O}(1)$ amortized)

 $\rightsquigarrow \mathcal{O}(m + n \log n)$ using Fibonacci Heaps

Problem: inherently sequential.

Best bet: use $\log n$ procs to support $\mathcal{O}(1)$ time PQ access.



Kruskal's Algorithm [1956]



 $T := \emptyset$

// subforest of the MST

foreach $(u, v) \in E$ in ascending order of weight **do**

if u and v are in different subtrees of T then

 $T := T \cup \{(u, v)\}$ // Join two subtrees

return T

Analysis



 $\mathscr{O}(\operatorname{sort}(m) + m\alpha(m, n)) = \mathscr{O}(m\log m)$ where α is the inverse Ackermann function

Problem: still sequential

Best bet: parallelize sorting

Idea: grow tree more aggressively



Edge Contraction

Let $\{u, v\}$ denote an MST edge. Eliminate *v*:

forall $(w, v) \in E$ do $E := E \setminus (w, v) \cup \{(w, u)\}$ // but remember orignal terminals





Boruvka's Algorithm

[Boruvka 26, Sollin 65]

For each node find the lightest incident edge. Include them into the MST (cut property) contract these edges,

Time $\mathcal{O}(m)$ per iteration

At least halves the number of remaining nodes





Analysis (Sequential)

 $\mathcal{O}(m\log n)$ time

asymptotics is OK for sparse graphs

Goal: $\mathcal{O}(m \log n)$ work $\mathcal{O}(\text{Polylog}(m))$ time parallelization



Finding lightest incident edges

Assume the input is given in adjacency array representation

```
forall v \in V dopar
```

allocate $|\Gamma(v)| \frac{p}{2m}$ processors to node v // prefix sum find w such that c(v, w) is minimized among $\Gamma(v)$ // reduction output original edge corresponding to (v, w)pred(v) := w

Time
$$\mathscr{O}\left(\frac{m}{p} + \log p\right)$$





Structure of Resulting Components

Consider a component *C* of the graph $(V, \{(v, pred(v)) : v \in V\})$

out-degree 1

 $\Box |C|$ edges

] pseudotree,

i.e. a tree plus one edge

one two-cycle at the lightest edge (u, w)

remaining edges lead to *u* or *w*





$\textbf{Pseudotrees} \rightarrow \textbf{Rooted Trees}$

forall $v \in V$ dopar

 $w := \operatorname{pred}(v)$ if $v < w \land \operatorname{pred}(w) = v$ then $\operatorname{pred}(v) := v$









Rooted Trees \rightarrow **Rooted Stars** by **Doubling**

while $\exists v \in V : \operatorname{pred}(\operatorname{pred}(v)) \neq \operatorname{pred}(v)$ do forall $v \in V$ dopar $\operatorname{pred}(v) := \operatorname{pred}(\operatorname{pred}(v))$





Contraction

k := # componentsV' = 1..k

Time $\mathscr{O}\left(\frac{m}{p} + \log p\right)$

$$\begin{split} \text{find a bijective mapping } f: \text{star-roots} &\to 1..k \\ E' &\coloneqq \{(f(\text{pred}(u)), f(\text{pred}(v)), c, e_{\text{old}}): \\ (u, v, c, e_{\text{old}}) \in E \land \text{pred}(u) \neq \text{pred}(v)\} \end{split}$$





Recursion



convert G' = (V', E') into adjacency array representation// integer sorting optional: remove parallel edges // retain lightest one recurse on G'

Expected sorting time $\mathscr{O}\left(\frac{m}{p} + \log p\right)$ CRCW PRAM [Rajasekaran and Reif 1989]

practical algorithms for $m \gg p$



Analysis

Satz 5. On a CRCW-PRAM, parallel Borůvka can be implemented to run in expected time

$$\mathscr{O}\left(\frac{m}{p}\log n + \log^2 n\right)$$

٠

$\Box \leq \log n$ iterations



☐ For root finding:

$$\sum_{i} \frac{n}{2^{i}} \log \frac{n}{2^{i}} \le n \log n \sum_{i} 2^{-i} = \mathscr{O}(n \log n)$$



Randomized Linear Time Algorithm

- 1. Factor 8 node reduction (3× Boruvka or sweep algorithm) $\mathscr{O}(m+n)$.
- 2. $R \leftarrow m/2$ random edges. $\mathcal{O}(m+n)$.
- 3. $F \Leftarrow MST(R)$ [Recursively].
- 4. Find light edges *L* (edge reduction). $\mathscr{O}(m+n)$ $\mathrm{E}[|L|] \leq \frac{mn/8}{m/2} = n/4.$

5. $T \leftarrow MST(L \cup F)$ [Recursively].

 $T(n,m) \leq T(n/8,m/2) + T(n/8,n/4) + c(n+m)$ $T(n,m) \leq 2c(n+m)$ fulfills this recurrence.

Parallel Filter Kruskal

Procedure filterKruskal(E, T : Sequence of Edge, P : UnionFind)

if $m \leq kruskalThreshold(n, m, |T|)$ then

kruskal(E,T,P)

else

pick a pivot $p \in E$ $E_{\leq} := \langle e \in E : e \leq p \rangle$ $E_{>} := \langle e \in E : e > p \rangle$ filterKruskal (E_{\leq}, T, P) **if** |T| = n - 1 **then** exit $E_{>} := \text{filter}(E_{>}, P)$ filterKruskal $(E_{>}, T, P)$

// parallel
// partitioning

// parallel sort

// parallel removeIf





Running Time: Random graph with 2¹⁶ **nodes**



number of edges m / number of nodes n

280



More on Parallel MST

[Pettie Ramachandran 02] $\mathcal{O}(m)$ work, $\mathcal{O}(\log n)$ expected time randomized EREW PRAM algorithm.

[Masterarbeit Wei Zhou 17:] Parallel Borůvka + filtering.

Use edge list representation, union-find.

Speedup up to 20 on 72 cores.

Master thesis topic: communication-efficient distributed-memory MST



Load Balancing

[Sanders Worsch 97]

Given

Work to be done

PEs

Load balancing = assigning work \rightarrow PEs

Goal: minimize parallel execution time

What we Already saw

- Estimating load using sampling
- Assign approx. equal sized pieces

sample sort

sample sort

Multi sequence selection balances multiway merging

Dynamic load balancing for quicksort and doall

Prefix sums

quicksort, parPQ, list ranking, MSTs,...

Parallel priority queues

branch-and-bound





Measuring Cost

$$\square \text{ Maximal Load:} \max_{i=1}^{p} \sum_{j \in \text{jobs @ PE } i} T(j, i, ...)$$

Time for finding the assignment

Executing the assignment

Cost for redistribution

communication between jobs? (volume, locality?)



What is Known About the Jobs?

Exact size

Approximate size

(Almost) nothing

Further subdivisible?

Similar for communication costs



What is Kown About the **Processors**?

All equal?

Different?

- Fluctuating external load
- Tolerate faults?

Similar for communication resources



In this Lecture

Independent jobs

- Sizes exactly kown fully parallel implementation
- Sizes unknown or inaccurate random assignment, master-worker-scheme, random polling



A very simple modell

- □ *n* jobs, $\mathcal{O}(n/p)$ per PE, independent, splittable, description with size $\mathcal{O}(1)$
- \Box Size ℓ_i exactly known


Sequential Next Fit [McNaughton 59]

$C := \sum_{j \le n} \frac{\ell_i}{p}$	// work per PE
i = 0	// current PE
f := C	// free room on PE <i>i</i>
j := 1	// current Job
$\ell := \ell_1$	// remaining piece of job j
while $j \le n$ do	
$c := \min(f, \ell)$	<pre>// largest fitting piece</pre>
assign a piece of size c of job j to ${\sf I}$	PE i
f := f - c	
$\ell := \ell - c$	
if $f = 0$ then $i + +; f := C$	// next PE
if $\ell = 0$ then $j + +$; $\ell := \ell_j$	// next job

289

Sanders: Parallel Algorithms January 17, 2022



Sequential Next Fit [McNaughton 59]





Parallelization of Next Fit (Sketch)

// Assume PE i holds jobs j_i..j'_i C:= $\sum_{j \leq n} \frac{\ell_i}{p}$ forall $j \leq n$ dopar pos:= $\sum_{k < i} \ell_k$ // prefix sums assign job j to PEs $\lfloor \frac{\text{pos}}{C} \rfloor$.. $\lfloor \frac{\text{pos} + \ell_j}{C} \rfloor$ // segmented broadcast piece size at PE $i = \lfloor \frac{\text{pos}}{C} \rfloor$: (i+1)C - pos piece size at PE $i = \lfloor \frac{\text{pos} + \ell_j}{C} \rfloor$: pos + $\ell_j - iC$

Time $C + \mathcal{O}\left(\frac{n}{p} + \log p\right)$ if jobs are initially distributed randomly.



292

Parallelization of Next Fit: Example





Atomic Jobs

Assign job j to PE $\left\lfloor \frac{\text{pos}}{C} \right\rfloor$

Maximum load $\leq C + \max_{j} \ell_{j} \leq 2$ opt

Better sequential approximation:

Assign largest jobs first

(shortest queue, first fit, best fit) in time $\mathcal{O}(n \log n)$

probably not parallelizable

Parallel

$$\frac{11}{9} \cdot \text{opt}$$

[Anderson, Mayr, Warmuth 89]



Atomic Jobs: Example



294



Example Mandelbrot Set

$$z_c(m) : \mathbb{N} \to \mathbb{C}$$

$$z_c(0) := 0, \quad z_c(m+1) := z_c(m)^2 + c$$

$$M := \{c \in \mathbb{C} : z_c(m) \text{ is bounded}\}.$$





296



Computation only for a quadratic subset of the complex plane

- Computation only for a discrete grid of points
- $\exists z_c$ unbounded if $|z_c(k)| \ge 2$
- Stop after m_{max} iterations

Where is the load balancing problem?





Code

```
int iterate(int pos, int resolution, double step)
{ int iter;
  complex c =
    z0+complex((double)(pos % resolution) * step,
                (double) (pos / resolution) * step);
  complex z = c;
  for (iter = 1;
       iter < maxiter && abs(z) <= LARGE;</pre>
       iter++) {
    z = z \star z + c;
  }
  return iter; }
```

Static Apple Distribution

298

Since there is little communication, we are very flexible

Split into strips

- Why attractive?
- Why rather not?
- Cyclic. Good. But provable??

```
Random
```



=PE 0









Parallelize Assignment Phase

- If the jobs are arbitrarily distributed over the PEs: Random permutation via all-to-all. (see also sample sort)
- Implicit generation of jobs
 - Jobs can be generated based on a number $1 \dots n$.
 - Problem: Parallel generation of a (pseudo) random permutation



Wlog (?) let n be a square.

Interpret numbers from 0..n-1 as pairs from $\{0..\sqrt{n}-1\}^2$.

 $\Box f: 0..\sqrt{n-1} \rightarrow 0..\sqrt{n-1}$ (pseudo)random function

Feistel permutation: $\pi_f((a,b)) = (b, a + f(b) \mod \sqrt{n})$ $(\pi_f^{-1}(b,x) = (x - f(b) \mod \sqrt{n}, b))$

Chain several Feistel permutations

$$\square \ \pi(x) = \pi_f(\pi_g(\pi_h(\pi_l(x))))$$
is even safe in some cryptographical sense



Random Assignment

Given: *n* jobs of size
$$\ell_1, \ldots, \ell_n$$

 $\Box \text{ Let } L := \sum_{i \le n} \ell_i$ $\Box \text{ Let } l_{\max} := \max_{i \le n} \ell_i$

Assign the jobs to random PEs

Theorem: If $L \ge 2(\beta + 1)pl_{\max}\frac{\ln p}{\epsilon^2} + O(1/\epsilon^3))$ then the maximum Load is at most $(1 + \epsilon)\frac{L}{p}$ with probability at least $1 - p^{-\beta}$. Proof: ... Chernoff-Bounds... Sanders: Parallel Algorithms January 17, 2022

Discussion



- + Job sizes need not been known at all
- + It is irrelevant where the jobs come from (distributed generation possible)
- Inacceptable with big l_{max}
- Very good load balance only with large L/l_{max} (quadratic in $1/\epsilon$, logarithmic in p).



Application Example: Airline Crew Scheduling

A single random assignment solves k simultaneous load balancing problems. (Deterministically, this is probably a difficult problem.)



Sanders: Parallel Algorithms January 17, 2022

The Master-Worker-Scheme





Discussion



- + Simple
- + Natural input-output scheme (but perhaps a separate disk slave)
- + Suggests itself when the job generator is not parallelized
- + Easy to debug
- Communication bottleneck \Rightarrow tradeoff communication cost versus imbalance
- How to split?
- Multi-level schemes are complicated and of limited help



Size of the jobs

- Deal out jobs that are as large as possible as long as balance is not in danger. Why?
- Conservative criterion: upper bound for the size of the delivered jobs \leq 1/P-th part of lower bound
 - for system load.
- Where to get size estimate?
- More aggressive approaches make sense





Work Stealing

- ☐ (Almost) arbitrarily subdivisible load
- Initially all the work on PE 0
- Almost nothing is known on job sizes
- Preemption is allowed. (Successive splitting)

Sanders: Parallel Algorithms January 17, 2022



Example: The 15-Puzzle



Korf 85: Iterative deepening depth first search with $\approx 10^9$ tree nodes.

Example: Firing Squad Synchronization Problem

#G.#G..#G...#G...#G....#G....#G....# #GG#GX.#GX..#GX...#GX...#GX....#GX....# #FF#XXX#XXG.#XXG..#XXG...#XXG....#XXG....# #GGG#GX.G#GX.X.#GX.X..#GX.X...#GX.X...# #FFF#XXXX#XXG.X#XXG.G.#XXG.G..#XXG.G...# #GGGG#GX.GX#GX.XXG#GX.XXX.#GX.XXX.# #FFFF#XXXX#XXG.XX#XXG.G.X#XXG.G.G.# #GGGGG#GX.G.G#GX.XXGX#GX.XXXXG# #FFFFF#XXXXX#XXG.XGX#XXG.GGXX# #GGGGGG#GX.G.GX#GX.XXGXG# #FFFFF#XXXXXX#XXG.XGXX# #GGGGGGGG#GX.G.GXG# #FFFFFF#XXXXXXX# #GGGGGGGG# # ㅋㅋㅋㅋㅋㅋㅋ

Backtracking over Transition Functions





Goal for the analysis

$$T_{\mathrm{par}} \leq (1 + \varepsilon) \frac{T_{\mathrm{seq}}}{p} + \mathrm{lower \ order \ terms}$$



An Abstract Model: Tree Shaped Computations





Tree Shaped Computations: Parameters

 T_{atomic} : max. time for finishing up an atomic subproblem

 T_{split} : max. time needed for splitting

h: max. generation gen(P) of a nonatomic subproblem *P*

 ℓ : max size of a subproblem description

p: no. of processors

 $T_{\rm rout}$: time needed for communicating a subproblem ($lpha + \ell eta$)

 T_{coll} : time for a reduction



Relation to Depth First Search

let stack consist of root node only

while stack is not empty do

remove a node \boldsymbol{N} from the stack

 $\mathbf{if} \ N \ \mathbf{is} \ \mathbf{a} \ \mathbf{leaf} \ \mathbf{then}$

evaluate leaf N

else

put successors of N on the stack

fi



Splitting Stacks





Other Problem Categories

Loop Scheduling

Higher Dimensional Interval Subdivision

Particle Physics Simulation

 \Box Generalization: Multithreaded computations. $h \rightsquigarrow T_{\infty}$

An Application List

- Discrete Mathematics (Toys?):
 - Golomb Rulers
 - Cellular Automata, Trellis Automata
 - 15-Puzzle, *n*-Queens, Pentominoes ...
- NP-complete Problems (nondeterminism ~> branching)
 - 0/1 Knapsack Problem (fast!)
 - Quadratic Assignment Problem
 - SAT
 - Functional, Logical Programming Languages
 - Constraint Satisfaction, Planning, ...
- Numerical: Adaptive Integration, Nonlinear Optimization by Interval Arithmetics, Eigenvalues of Tridiagonal Matrices





Limits of the Model

- Quicksort and similar divide-and-conquer algorithms (shared memory OK → Cilk, MCSTL, Intel TBB, OpenMP 3.0?)
 - Finding the first Solution (often OK)
 - Branch-and-bound
 - Verifying bounds OK
 - Depth-first often OK
- Subtree dependent pruning
 - FSSP OK
 - Game tree search tough (load balancing OK)



Receiver Initiated Load Balancing





Random Polling





$\tilde{O}\left(\cdot\right)$ Calculus

 $X \in \tilde{\mathrm{O}}\left(f(n)
ight)$ – iff $orall oldsymbol{\beta} > 0$:

$$\exists c > 0, n_0 > 0 : \forall n \ge n_0 : \mathbb{P}\left[X > cf(n)\right] \le n^{-\beta}$$

Advantage: simple rules for sum and maximum.

Sanders: Parallel Algorithms January 17, 2022



not here





Synchronous Random Polling

P, P': Subproblem $P := \text{if } i_{\text{PE}} = 0 \text{ then } P_{\text{root}} \text{ else } P_{\emptyset}$ $P := \operatorname{work}(P, \Delta t)$ loop $m' := |\{i : T(P@i) = 0\}|$ if m' = p then exit loop else if m' > m then if T(P) = 0 then send a request to a random PE if there is an incoming request then $(P, P') := \operatorname{split}(P)$ send P' to one of the requestors send empty subproblems the rest if T(P) = 0 then receive P



Analysis

Theorem:

For all $\varepsilon > 0$ there is a choice of Δt and m such that

$$T_{\text{par}} \leq (1+\varepsilon) \frac{T_{\text{seq}}}{p} + \tilde{O}\left(T_{\text{atomic}} + h(T_{\text{rout}}(l) + T_{\text{coll}} + T_{\text{split}})\right)$$




Bounding Idleness

Lemma 6.

Let m < p with $m \in \Omega(p)$. Then $\tilde{O}(h)$ iterations with at least m empty subproblems suffice to ensure $\forall P : gen(P) \ge h$.





Busy phases

Lemma 7. There are at most $\frac{T_{\text{seq}}}{(p-m)\Delta t}$ iterations with $\leq m$ idle *PEs at their end.*



A Simplified Algorithm

P, P': Subproblem

 $P := \mathbf{if} \ i_{\text{PE}} = 0 \mathbf{then} \ P_{\text{root}} \mathbf{else} \ P_{\emptyset}$

while not finished

 $P := \operatorname{work}(P, \Delta t)$

select a global value $0 \le s < n$ uniformly at random if $T(P@i_{PE} - s \mod p) = 0$ then $(P, P@i_{PE} - s \mod p) := \operatorname{split}(P)$

Satz 8. For all $\varepsilon > 0$ there is a choice of Δt and m such that

$$T_{\text{par}} \preceq (1+\varepsilon) \frac{T_{\text{seq}}}{p} + \tilde{O}\left(T_{\text{atomic}} + h(T_{\text{rout}}(l) + T_{\text{split}})\right)$$



328

Asynchronous Random Polling

P, P': Subproblem

 $P := \mathbf{if} \ i_{\text{PE}} = 0 \mathbf{then} \ P_{\text{root}} \mathbf{else} \ P_{\emptyset}$

while no global termination yet $d\boldsymbol{o}$

if T(P) = 0 then send a request to a random PE

else $P := \operatorname{work}(P, \Delta t)$

 ${f if}$ there is an incoming message M then

if *M* is a request from PE j **then** $(P, P') \coloneqq \operatorname{split}(P)$ send *P'* to PE j

else

$$P := M$$



Analysis

Satz 9.

$$\mathbb{E}T_{\text{par}} \leq (1+\varepsilon)\frac{T_{\text{seq}}}{p} + \mathcal{O}\left(T_{\text{atomic}} + h\left(\frac{1}{\varepsilon} + T_{\text{rout}} + T_{\text{split}}\right)\right)$$

for an appropriate choice of Δt .



A Trivial Lower Bound

Satz 10. For all tree shaped computations

$$T_{\text{par}} \in \Omega\left(\frac{T_{\text{seq}}}{p} + T_{\text{atomic}} + T_{\text{coll}} + T_{\text{split}}\log p\right)$$
.

if efficiency in $\Omega(1)$ *shall be achieved.*



Many Consecutive Splits



Additional

$$h - \log \frac{T_{\text{seq}}}{T_{\text{atomic}}}$$

term.



Many Splits Overall





Satz 11. Some problems need at least

$$\frac{p}{2} \left(h - \log \frac{T_{\text{seq}}}{T_{\text{atomic}}} \right)$$

splits for efficiency $\geq \frac{1}{2}$.

Korollar 12. *Receiver initiated algorithms need a corresponding number of communications.*

Satz 13 (Wu and Kung 1991). *A similar bound holds for all deterministic load balancers*.



Golomb-Rulers



Applications: Radar astronomy, codes, ...



Many Processors

Parsytec GCel-3/1024 with COSY (PB)

Verification search



LAN



Differing PE-Speeds (even dynamically) are unproblematic.

Even complete suspension OK as long as requests are answered.





The 0/1**-Knapsack Problem**

] *m* items

] Maximum knapsack weight M

Item weights w_i

 \Box Item profits p_i

- \Box Find $x_i \in \{0, 1\}$ such that
 - $-\sum w_i x_i \leq M$
 - $\sum p_i x_i$ is maximized



Best known approach for large *m*:

- Depth-first branch-and-bound
- Bounding function based on a the relaxation $x_i \in [0, 1]$. (Can be computed in $\mathcal{O}(\log m)$ steps.)



Superlinear Speedup

Parsytec GCel-3/1024 under COSY	(PB)
---------------------------------	------

- 1024 processors
- 2000 items
- Splitting on all levels
- 256 random instances at the border between simple and difficult
- Overall $1410 \times$ faster than seq. computation!







Fast Initialization





Static vs Dynamic LB





Beyond Global Polling



Asynchronously increase polling range (exponentially)





Scalability Comparison Independ. Jobs

- \overline{t} =average job size
- $\hat{t} =$ maximum job size

T =required total work for parallel execution time $(1 + \varepsilon) \frac{\text{total work}}{p}$

Algorithm	$T = \Omega(\cdots)$	Remarks
prefix sum	$\frac{p}{\epsilon}(\hat{t} + \alpha \log p)$	known task sizes
master-worker	$\frac{p}{\varepsilon} \cdot \frac{\alpha p}{\varepsilon} \cdot \frac{\hat{t}}{\bar{t}}$	bundle size $\sqrt{\frac{m\alpha}{\hat{t}}}$
randomized static	$\frac{p}{\varepsilon} \cdot \frac{\log p}{\varepsilon} \cdot \hat{t}$	randomized
work stealing	$\frac{p}{\varepsilon}(\hat{t} + \alpha \log p)$	randomized





-] Naive Parallelization yields only Speedup $\mathscr{O}(\sqrt{n})$.
- Young Brother Wait Concept (Feldmann et al.)
- Tradeoff between Speculativity and Sequentialization
- Propagate window updates
- Combine with global transposition table



MapReduce in 10 Minutes

[Google, DeanGhemawat OSDI 2004] see Wikipedia

Framework for processing multisets of (key, value) Pairs.



 $// M \subset K \times V$ // MapF : $K \times V \rightarrow K' \times V'$ // ReduceF : $K' \times 2^{V'} \rightarrow V''$ **Function** mapReduce(M, MapF, ReduceF) : V'' $M' := \{ MapF((k,v)) : (k,v) \in M \}$ // easy (load balancing?) $\operatorname{sort}(M')$ // basic toolbox forall k' with $\exists (k', v') \in M'$ dopar *ll* easy $s := \{v' : (k', v') \in M'\}$ $S := S \cup (k', s)$ **return** {reduceF $(k', s) : (k', s) \in S$ } // easy (load balancing?)



Refinements

Fault Tolerance

Load Balancing using hashing (default) und Master-Worker

Associative commutative reduce functions

Examples



- URL access frequencies
- build inverted index
- Build reverse graph adjacency array





Graph Partitioning

Contraction

while $|V| > c \cdot k$ do find a matching $M \subseteq E$ contract M // similar to MST algorithm (more simple)

save each generated level



Finding a Matching

Find approximate max. weight matching wrt edge rating

$$expansion(\{u,v\}) \coloneqq \frac{\omega(\{u,v\})}{c(u) + c(v)}$$

$$expansion^*(\{u,v\}) \coloneqq \frac{\omega(\{u,v\})}{c(u)c(v)}$$

$$expansion^{*2}(\{u,v\}) \coloneqq \frac{\omega(\{u,v\})^2}{c(u)c(v)}$$

$$innerOuter(\{u,v\}) \coloneqq \frac{\omega(\{u,v\})}{Out(v) + Out(u) - 2\omega(u,v)}$$



Approx. max. weighted Matching

todo



Graph Partitioning Future Work

- Understand edge ratings
- Scalable parallel weighted Matching code
- Hypergraph partitioning
- Handling exact balance
- Max. Flow. based techniques
- Parallel external, e.g., partitioning THE web graph