

Zaawansowane Techniki Algorytmiczne

Mirosław Kutylowski 2017

Katedra Informatyki WPPT, PWr

Plan wykładu

1. modele obliczeń
 - a. systemy równoległe
 - b. systemy rozproszone
 - c. sieci Boolowskie, OBDD
 - d. obliczenia kwantowe i inne modele odbiegające od modelu von Neumanna
2. paradygmaty algorytmiczne
 - a. samostabilizacja
 - b. algorytmy aproksymacyjne
 - c. algorytmy dla danych rozmytych
 - d. algorytmy losowe
 - e. derandomizacja
 - f. algorytmy online
 - g. uniwersalne heurystyki
3. analiza złożoności
 - a. granice dolne
 - b. złożoność komunikacyjna
 - c. kombinatoryka analityczna
 - d. rapid mixing
 - e. rozwiązania dla ograniczonych zasobów

Cele

wykształcenie umiejętności w zakresie wykorzystania szerokiego spektrum zaawansowanych technik algorytmicznych, umiejętności z zakresu analizy poprawności i efektywności algorytmów

MODELS

I. Parallel computing

parallel computing today:

- supercomputers
- clusters of computers tightly connected (supercomputing centers)
- multicore architectures
- hardware like graphic cards
- some embedded systems

limitations for increasing computing power:

- each operation consumes energy, physical limitations
- increasing speed (frequency of the processor clock) increases energy consumption per second on a mm^2
- the same effect of more dense layout - large scale of integration
- ... but the heat has to be removed in order to prevent overheating. This is hard (cooling systems)!
- \Rightarrow progress on single processor architectures has been stopped after rapid progress during the 90's
- parallelisation as a remaining option

types of parallel machines:

- program execution:
 - SIMD - single instruction multiple data (all processors execute the same code and are in the same place in the program, typically used for „vector computations“, that is matrices and so)
 - MIMD - multiple instruction multiple data
- memory type:
 - shared (all locations accessible by all processors, problems of conflicts)
 - distributed (each processor has own memory, addressing remote memory via its owner)
 - hybrid (both types, shared memory slower and problematic, but sometimes extremely useful)

- interconnections:
 - mesh
 - hypercube or its variants (butterfly, deBruijn)
 - fat tree
- communication between the processors:
 - MPI (Message Passing Interface) a standard for point-to-point communication, buffers on endpoints, coordination: no message delivered before it is sent
 - shared memory (types: concurrent writes, collision, exclusive write,...)
- coordination between processors:
 - MPI: the programmer is fully responsible for it, logical structure of the algorithm has to ensure proper execution
 - shared memory: coordination from the global clock but ... latency, conflicts, ...
 - example: computing OR of n bits on a PRAM machine:
 - in time t each processor knows OR of $P(t)$ values
 - in time t each memory cell stores OR of $M(t)$ values
 - after writing: $M(t+1) = P(t) + M(t)$
 - after reading: $P(t+1) = M(t+1) + P(t)$
 - a Fibonacci sequence ... $M(t), P(t), M(t+1), P(t+1)$
 - grows almost like exponentially, but not exactly
 - so a careful design which address to read and write
 - BSP (Bulk Synchronous Processing): supersteps, at the end of the superstep everything synchronised again, during the superstep the timing unpredictable
 - LogP: a model with synchronous clock and limitation of latency of communication
 - L communication latency (time to deliver a message)
 - o - overhead (time needed to send or receive a message), at this time no other operation performed
 - g - gap, time between consecutive transmissions of messages
 - Example: computing $a_1 \times a_2 \times \dots \times a_n$ for an associative operation \times :
 - obviously, a tree structure for collecting data and combining results

- we determine how big might be n given a time limit T :
 - for $T \leq L + 2o$ communication would cost more, so compute everything locally
 - otherwise the last step of the root is to combine its result and the result obtained from another processor, which has sent at time: $T - 1 - L - 2o$
 - ... so we may get a recursive expression for the maximal n

- parallelisation:
 - by hand: define processes (processes are assigned to processors by the manager of a parallel machine)
 - writing a code, giving to a compiler that recognizes parallelism and transfers into appropriate code
 - necessary to write programs with explicit parallelism in mind
 - in some languages explicit declarations
 generally: **a weak point (people think sequentially)**

application areas:

- numerical computations of linear algebra
- all spacial computations (2D, 3D,..) for instance in civil engineering, weather forecast,
- crypto - code breaking via brute force and many search algorithms
- bioinformatics, genetics, chemistry (simulating as a cheap and fast initial stage of design to eliminate most of wrong directions)
- modelling complex systems

computational complexity for parallel computing:

- **time** (sometimes called *depth*) from the start to the moment when the output is ready
 - if there is no fixed termination time it might be problematic to reach consensus when the computation has been finished
- **work** = the total amount of steps executed by all processors involved
 - generally the work cannot be lower than in case of a sequential execution (a sequential machine can always simulate a parallel computation)

- typically the price for time speed-up is an increase of work

main problems:

- time to market
- correctness of algorithms (practically infinite number of options for timing options of interprocessor communication)
- most programmers cannot think in terms of parallel programs
- **some problems cannot be solved faster by parallel machines**

Example of inherently non-parallel task:

compute H_n for a given I , t and n , where $H_0 = I$, $H_{i+1} = \text{Hash}(H_i, i, t)$
 Hash is a cryptographic hash function

Knowledge assumption (simplified version, one of fundamental assumptions for many cryptographic constructions):

if a system learns h and y such that $h = \text{Hash}(y)$, then with very high probability y must appear during the computation **before** h

Many algorithmic problems seem to be hard to parallelize. It is an open theoretical question whether, say, all problems solvable in polynomial time can be executed in much shorter time by a parallel machine.

Some issues:

- **load balancing** (utilize processors evenly):
 - sophisticated approaches or
 - random strategies (e.g. choose two processors at random and assign the task to a processor with less load – so called „power of two choices“)
- **symmetry breaking:** if no IDs assigned to processes then processes might be in exactly the same state and who should take which role. Solution: based on random choice
- **consensus:** processes might have a different understanding of the global state
- **Byzantine agreement:** the messages might be undelivered. What to do in this case? A majority 2/3 of nonmalicious processors should be able to make a decision

(Byzantine agreement problem: there are two Byzantine armies. If they attack the enemy at the same time, then they win. Communication between the armies is via messengers that go through the enemy territory. The messengers can be captured - the sender cannot be sure that the message arrived at the destination)

ALGORITHM EXAMPLES

(borrowed from *Parallel Algorithms* by G. Blelloch and B. Maggs)

PARALLEL PREFIX: compute $\sum_{i=1}^j A[i]$, simultaneously for $j = 1, 2, \dots, n$

ALGORITHM: ParallelPrefix(A)

- 1 if $|A| = 1$ then return $A[0]$
- 2 else
- 3 $S = \text{ParallelPrefix}(\{A[2i] + A[2i + 1]\}_{i=1, \dots, |A|/2})$
- 4 $R[i] := S[i/2]$ if i even, else $R[i] := S[i - 1/2] + A[i]$ for $i \leq n$

array R is the output

Work: $W(n) = W(n/2) + O(n)$ so $W(n) = O(n)$

Time $D(n) = D(n/2) + O(1)$ so $D(n) = O(\log n)$

POINTER-JUMPING: given a directed acyclic graph (e.g. a tree), find the root for each vertex

ALGORITHM: pointer-jumping(P)

- 1 for j from 1 to $\lceil \log |P| \rceil$
- 2 $P := P[P[i]]$ for $i \leq |P|$

- at each iteration the pointer jumps forwards, the exception are the roots that point to themselves
- if the pointer already is to the root then the pointer does not change
- generally, the length of jumps double at each iteration, as the maximal path has length $|P|$, no more than $\log |P|$ iterations are needed

LIST-RANKING: given a list represented via pointers, find the distance of each vertex from the head of the list

idea: like pointer jumping, but keep counting the distance to the node shown by the pointer

ALGORITHM: list-ranking(P)

- 1 assign $V[i] = 1$ unless $P[i] = i$ (pointer to itself)
- 2 for j from 1 to $\lceil \log |P| \rceil$
- 3 $V[i] := V[i] + V[P[i]]$ for $i \leq |P|$

4 $P := P[P[i]]$ for $i \leq |P|$

V is the output

Work is $\Theta(n \log n)$. Bad!

a typical improvement via random sampling:

- choose $n/\log n$ start nodes at random
- from each start node walk (1 process per start node, the walk is sequential) until another start node is encountered
- with high probability each walk stops after $O(\log n)$ steps
- solve the LIST-RANKING problem with the list of start nodes and initialized not with 1 but the distance to the next start node
- for the reduced problem the previous algorithm can be applied, it requires $O(n/\log n \cdot \log(n/\log n)) = O(n)$ work, – so within the bound $O(n)$
- walk back (in parallel starting from each starting node) and compute the distances on the way

execution time remains logarithmic, but the work is time·number of processors = $O(\log n \cdot n/\log n) = O(n)$

REMOVING DUPLICATES: array contains entries, some of them appear more than once. Remove duplicates (leave only one position for a given value)

- if the range of the elements is small, the problem is easy to solve:
 - (in parallel) read a position in the input array,
 - if z found then write z into the output array $R[z] := z$concurrent write is necessary

a solution based on hashing:

ALGORITHM: remove-duplicates(V)

1 choose a prime m higher than $2 \cdot |V|$

2 fill TABLE with -1

3 $i := 0$

4 $R := \{\}$

5 while $|V| > 0$

6 in TABLE insert value j in position $\text{hash}(V[j], m, i)$ for each j

7 WINNERS := $\{V[j] : \text{TABLE}(\text{hash}(V[j], m, i)) = j\}$

8 append R with the list of winners

9 in TABLE insert $\text{hash}(k, m, i)$ in position k for each k from WINNERS

10 leave in V only those k , for which $\text{TABLE}[\text{hash}(k, m, i) \neq k]$
11 $i := i + 1$

RESULT is the output

- fine tuning regarding the choice of m at each stage (trade-off between efficiency - lower m means faster appending - and probability of collisions - low m means a lot of collisions and many rounds)
- appending the list of winners requires parallel prefix
- too small m means a lot of collisions
- each stage uses a different i , so hash values are unrelated

SORTING

- QUICKSORT is an inherently parallel algorithm
- but RADIXSORT works also fine (requires that the values are b bit numbers, time $O(b \cdot \log n)$, work $O(b \cdot n)$)

ALGORITHM: radixsort(A, b)

```
1 for  $i = 0, \dots, b - 1$ 
2 flags :=  $\{(a \gg i) \bmod 2 : a \in A\}$ 
3 notflags :=  $1 - \text{flags}$ 
4  $R_0 := \text{parallelprefix}(\text{notflags})$ 
5  $s_0 := \text{sum}(\text{notflags})$ 
6  $R_1 := \text{parallelprefix}(\text{flags})$ 
7  $R[j] := R_0[j]$  if  $\text{flags}[j] = 0$ , else  $R[j] := R_1[j] + s_0$  (computing ranks)
8 rewrite  $A$ : value  $A[j]$  moved to position  $R[j]$ 
```

A is the output

properties:

- stable sorting (order of the same elements preserved)
- first reorder according to the least significant bit, then 2nd, ...

CONNECTED-COMPONENTS:

definition: two vertices in a graph are in the same connected component if there is a path from one vertex to the second vertex

goal: label all vertices in a component with the same label, different label for different connected components

- one may try BFS search but it is inherently sequential and inefficient in the parallel setting

- it is easier to solve the problem using graph contraction technique

ALGORITHM: randomcontraction(LABELS, E)

```

1 if ( $|E| = 0$  then return LABELS
2 else
3  each vertex chooses a bit at random, a vertex called a child if 1 chosen
4  HOOKS :=  $\{(u, v) \in E: \text{child}[u] = 1 \text{ and } \text{child}[v] = 0\}$ 
5  put in LABELS values from HOOKS (the child gets the label of the parent)
6   $E' := \{(LABELS[u], LABELS[v]: (u, v) \in E \text{ and } LABELS[u] \neq LABELS[v]\}$ 
7  LABELS' := randomcontraction(LABELS,  $E'$ )
8  insert to LABELS' entries  $(u, LABELS'[v])$  for  $(u, v) \in \text{HOOKS}$ 

```

LABELS' returned

ALGORITHM: deterministiccontract(labels, E)

```

1 if ( $|E| = 0$  then return LABELS
2 else
3  HOOKS :=  $\{(u, v) \in E: u > v\}$ 
4  put in LABELS values from HOOKS (the child gets the label of the parent)
5  with pointer jumping assign LABELS according to the roots of local trees
6   $E' := \{(LABELS[u], LABELS[v]: (u, v) \in E \text{ and } \text{labels different}\}$ 
7  return deterministiccontract(LABELS,  $E'$ )

```

sometimes works badly:

- a star graph with labels $1..n-1$ outside and label n in the middle of the star
- n recursive calls are necessary

II. Distributed computing

(notices based on the script of R. Wattenhofer, ETH Zurich)

below 2 interesting general topics from distributed systems. These are only to examples.

Problems:

- communication takes long time

- communication is unreliable: messages may disappear or come in a wrong order
- network nodes may fail or even be malicious
- no direct coordination, the processors may have inconsistent view of the computation as a whole

Client-server Systems

- multiple users
- multiple servers – executing the commands of the users
- each command from a user has to be executed on all servers
- problem with the order of receiving commands:
 - servers X, Y store $a = 0$
 - client A asks to compute $a := a + 1$
 - client B asks to compute $a := 2 \cdot a$
 - X receives order from A and then from B , result: $a = 2$
 - Y receives order from B and then from A , result: $a = 1$
 - inconsistency created!!

Serializer:

- a distinguished server that collects the commands from the clients and send them to the servers in the same order
- problem: serializer is a single point of failure
- but the goal of deploying multiple servers was to avoid problem if some server fails!

Locks:

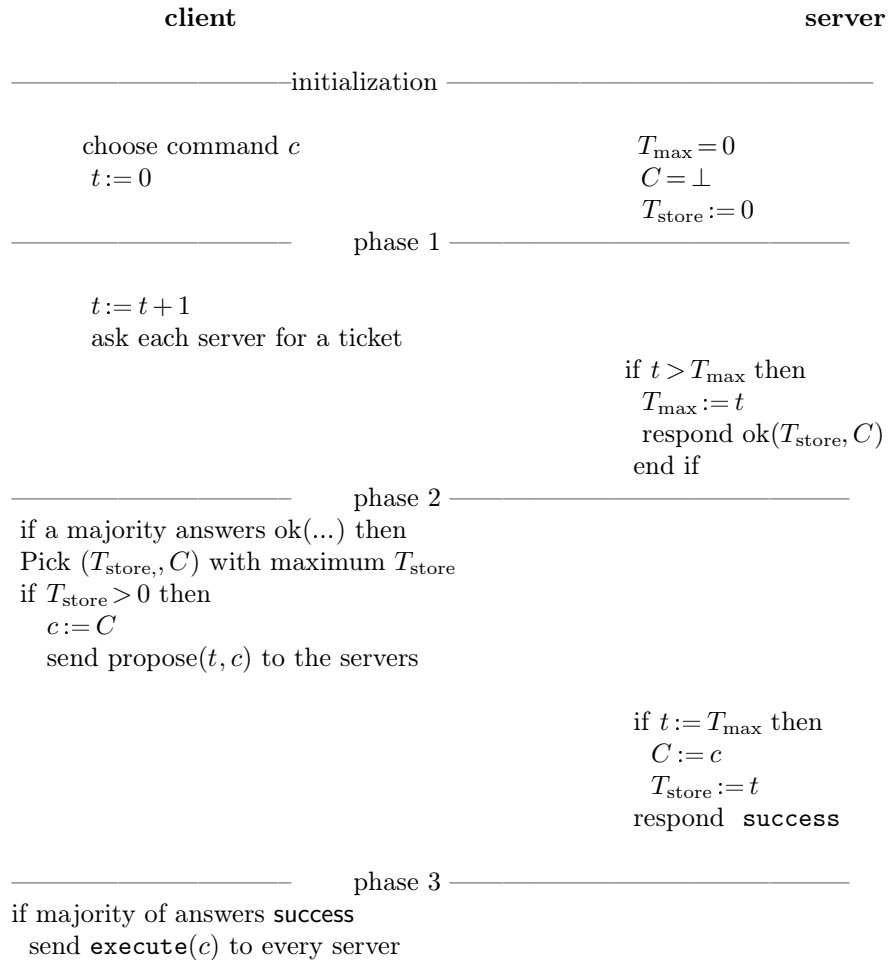
- each client tries first to lock all servers and then execute the command
- no inconsistency
- but what if some server does not respond? the system is blocked
- some strategy to unlock the servers needed (when two clients lock successfully different servers at the same time)

Idea of tickets:

- tickets issued to clients by the servers

- the tickets need not to be returned
- server accepts the commands with the most recently issued tickets only

PAXOS:



properties:

- if (t, c) accepted and stored by majority of servers, then later each **propose**(t', c') contains $c' = c$
 - take the smallest t^* for which this is not true
 - there must be server s that has been involved in both proposals
 - s received request for t^* after it has stored (t, c) (otherwise t^* would not be a valid ticket)

- but then the response would be $ok(t, c)$ and then $c' = c$
- no way to get $ok(t'', c')$ as t^* is the smallest one higher than t

Quorum Systems

Quorum system idea:

- there are n servers, data is to be stored on all of them
- however, doing it at once is hard
- if we update a record on a subset of servers (a *quorum*) - we lock it and update
- in order to read a data the user checks all servers from some *quorum* and takes the most recent one
- (we may assume that the data are authenticated together with time stamps so it is easy to see which version is the most recent one)
- different possibilities for quorums help as each server may fail to respond, ...

Problem:

- how to define quorums so that: any two quorums have non-empty intersection
- how to define access strategy: assign probability to each quorum and while reading choose a quorum according to this probability
- examples:
 - singleton: only one quorum consisting of one server
 - majority: every set of at least $n/2 + 1$ servers

Quality of a solution: **load**

- $L_Z(v)$ is load of a server v for a strategy Z : the probability that the server v will be read: $\sum_{v \in Q} \Pr(\text{strategy } Z \text{ chooses } Q)$
- $L_Z(S)$ – the load of a strategy Z over a quorum S : max over all $L_Z(v)$
- $L(S) = L_Z(S)$ for the best strategy Z

Quality of a solution: **work**

- $W_Z(S)$ – work for strategy Z is the expected size of a quorum chosen
- $W(S)$ – work for the quorum S for the best strategy Z

Examples:

singleton: work=1, load=1

majority: work $>n/2$, load ≈ 0.5

Theorem $L(S) \geq 1/\sqrt{n}$

- consider a quorum Q with the smallest size q
- claim: for some $v \in Q$, $L_Z(v) \geq 1/q$, indeed: each time a quorum is accessed at least one server from Q is accessed as well, the probabilities of quorums sum up to 1 and they are „distributed” among q servers
- each time at least q servers are accessed, so there must be a server with $L_Z(v) \geq q/n$
- $L_Z(v) \geq \max(1/q, q/n)$. The best choice is $q = \sqrt{n}$

Grid quorum system:

- the servers form a grid $d \times d$ ($n = d^2$)
- each quorum is a set consisting of a column and a row
- each two quorums have 2 points of intersection
- option 1: a row and a column truncated below this row (only one intersection guaranteed)
- option 2: one row plus server per row in the rows below
- load $\approx 2/\sqrt{n}$

Locking problem:

- each access has to lock the quorum before writing (otherwise a newer record might be overwritten by an older one)
- deadlock possible: e.g. in the grid system S_1 and S_2 intersect at s and s' :
 - S_1 locks s
 - S_2 locks s'
 - neither of them can proceed
- **Distributed Locking algorithm:**
 - lock the nodes of a quorum one by one, according to their id numbers in an increasing way
 - if a locked server encountered then release all servers locked so far
- Claim: no deadlock possible
Observation: the process that has locked the server with the highest id is either completed or can proceed (no node with a higher ID has been locked so far)

Fault tolerance

- up to f servers may fail, still there should be a quorum disjoint with the failed servers
- grid quorum system has f -resilience for $f < \sqrt{n}$, it is not \sqrt{n} fault resilient (fail nodes on the diagonal)

Probabilistic failure

- with pbb p a server is not in the failure state
- What is the probability that a quorum system S fails? Notation: $F_p(S)$
- behavior: $F_p(S)$ inspected for big n :
 - for majority quorum system $F_p(S) \rightarrow 0$ for $p < 0.5$
(follows from Chernoff Bound:
for m independent binary variables x_i and success pbb p)

$$\Pr\left(\sum_{i=1}^m x_i < (1 - \delta)m \cdot p\right) < e^{-m \cdot p \cdot \delta^2 / 2}$$

- for grid system, $F_p(S) \rightarrow 1$
 - we need at least one failed node per row to fail
 - $F_p(S) = (1 - p^d)^d > 1 - d \cdot p^d \rightarrow 1$ (as $(1 + x)^n > 1 + x \cdot n$ for $x > -1$)

B-Grid (a clever construction for low failure probability)

- a grid with $n = d \cdot h \cdot r$ nodes, d columns
- h bands, each consisting of r rows
- "minicolumn" is a column in a band
- quorum: a minicolumn in each band, additionally: in one chosen band: an element in each minicolumn
- $F_p(S) \rightarrow 0$
 - failure if in each band a complete minicolumn fails or in one band in each minicolumn an element fails:
$$F_p(S) \leq (d(1 - p)^r)^h + h(1 - p^r)^d$$
 - use $d = \sqrt{n}$, $r = \ln(d)$, $p \geq 2/3$
 - then $(d(1 - p)^r)^h \leq (d(1/3)^r)^h = d(d^{\ln 1/3})^h \approx d^{-0.1} < 1/n$
 - then $h(1 - p^r)^d < d(1 - p^r)^d < d(1 - d^{\ln 2/3})^d \approx d(1 - d^{-0.4})^d \approx d \cdot e^{-d^{0.6}} = d^{(-d^{0.6}/\ln d)+1} \ll d^{-2} \approx 1/n$

Byzantine systems

- up to f servers may cheat
- f -disseminating: every intersection of quorums contains at least $f + 1$ servers, there is a quorum without Byzantine nodes
- so always at least one quorum survives (the Byzantine ones may pretend to crash)
- after writing in one quorum and reading by another one there will be a witness of the correct value (the Byzantine nodes may store an old/wrong value)
- this is enough if data authenticated
- as in the proof of the theorem above:
$$L(S) \geq \sqrt{(f+1)/n}$$
- f -masking grid system: a quorum is a column and $f + 1$ rows, required:
$$2f + 1 \leq \sqrt{n}$$
- M-Grid: $\sqrt{f+1}$ rows and $\sqrt{f+1}$ columns, quorums intersection have
$$2\sqrt{f+1}^2 = 2(f+2)$$
 nodes

Opaque systems

- assume there is no data authentication
- two quorums intersect: Q_1 is up-to-date, Q_2 is not
- for any set F of f Byzantine nodes we require that the number of nodes in the intersection $Q_1 \cap Q_2 \setminus F$ is bigger than the number of nodes in $F \cap Q_2$ (cheating nodes) \cup $Q_2 \setminus Q_1$ (old values)
- would be great but...

Theorem. For any f opaque system $L(S) \geq 0.5$

Proof.

i. size of $Q_1 \cap Q_2$ is at least half of the size of Q_1 because of the condition on opaque systems

ii. load on Q_1 :

$$\sum_{v \in Q_1} \sum_{v \in Q_i} P_Z(Q_i) = \sum_i \sum_{v \in (Q_i \cap Q_1)} P_Z(Q_i) \geq \sum_i (|Q_1|/2) \cdot P_Z(Q_i) = |Q_1|/2$$

iii. now by pigeonhole principle: there must be some node in Q_1 with at least load 0.5