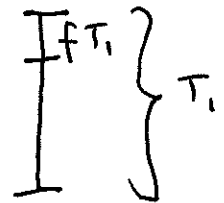


Reasoning about performance

Amdahl's Law (1967)

A program runs in time T_1 on one processor. A fraction f , $0 \leq f \leq 1$, of it is sequential. Let T_p be the runtime on p processors. Then



$$T_p \geq \frac{(1-f)T_1}{p} + fT_1$$

$$\text{speedup: } S_p = \frac{T_1}{T_p} \leq \frac{1}{\frac{1-f}{p} + f}$$

("linear" speedup means $S_p \approx p$)

$$\text{efficiency: } E_p = \frac{S_p}{p} \leq \frac{1}{1-f+fp}$$

($S_p \approx p \Rightarrow E_p \approx 1$)

$$p \rightarrow \infty: T_\infty \geq fT_1$$

$$S_\infty \leq \frac{1}{f}$$

$$E_\infty = 0 \quad \text{if } f \neq 0$$

Is Amdahl's Law optimistic or pessimistic?

Pessimistic

- AL fixes the problem size n . But with more processors people will want to solve larger problems. Thus take n into account above: $T_1(n), f(n), \dots$

For example (Gustafson's Law, 1988):

$$S_\infty(n) \leq \frac{1}{f(n)} \xrightarrow{n \rightarrow \infty} \infty \quad \text{if } f(n) \rightarrow 0$$

(sequential fraction decreases with problem size)

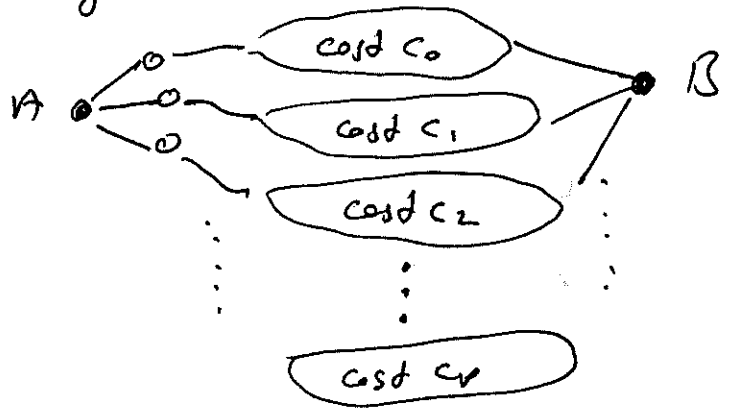
Terms:

- strong scaling: behavior of $S_p(n)$ for fixed n and $p \rightarrow \infty$
- weak scaling: " for $p, n \rightarrow \infty$

- AL assumes that by increasing p , all other resources remain constant. If this is not the case, super linear speedup is possible

- data caches private to a processor
(for $p=2$, working set can suddenly fit into cache)
- memory bandwidth
 $p=2$ threads on 2 processors may get twice the bandwidth

- AL does not allow algorithmic super linear speedup
- variant of traveling salesman: go from A to B visiting all cities



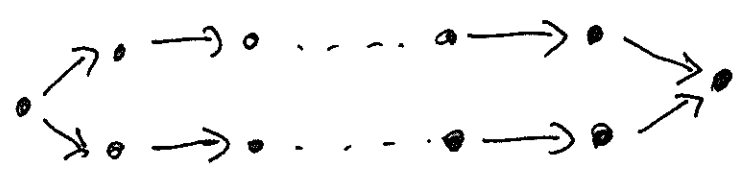
sequential: $\sum c_i$
 On p processors:
 $\min c_i$

Optimistic

- Ignores overhead of parallelization (e.g., creating threads) which increases with p .
- Assumes load balancing. In reality some processors wait.

So in reality: $S_p(n) = \frac{T_1(n)}{T_p(n) + A_p(n)}$ (overhead)

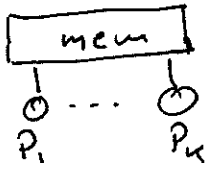
- Many programs have no sequential or infinitely parallelizable part. Example:



\Rightarrow Need better models that take graph structure into account.

PRAM machine model & DAG-based computation model

PRAM:



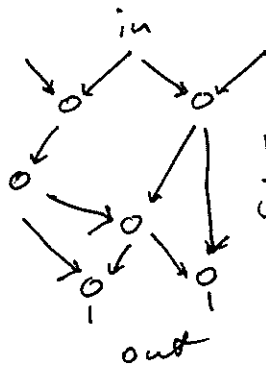
all processors can access every memory location simultaneously and in unit time

Computation as directed acyclic graph (DAG)

nodes: unit time operations

edges: dependencies

n : input size



$W(n) = 6$
 $D(n) = 4$

$W(n)$ (work, op count, cost):

number of nodes

$D(n)$ (depth, span, length of critical path)
longest path between any input and output (measured in # nodes on that path)

Average parallelism: $\frac{W(n)}{D(n)}$

Examples:

1.) Reduction: $x_0 + x_1 + \dots + x_{n-1}$

sequential: $W(n) = O(n), D(n) = O(n)$

binary tree: $W(n) = O(n), D(n) = O(\log(n))$

average parallelism

$O(1)$

$O(\frac{n}{\log n})$

2.) Mergesort: L a list of length n

sort(L)

if length(L) = 1 return L

$L_1 = \text{sort}(\text{left}(L))$

$L_2 = \text{sort}(\text{right}(L))$

return merge(L_1, L_2)

$W(n) = O(n \log(n))$

$D(n) = D(\frac{n}{2}) + O(n) = O(n)$

average parallelism:

$O(\log(n))$

Note: Parallel merge exists \rightarrow shorter $D(n) \rightarrow$ more parallelism

3.) Scan: input: list $x = (x_0, x_1, \dots, x_{n-1})$

output: $(0, x_0, x_0 + x_1, x_0 + x_1 + x_2, \dots, x_0 + \dots + x_{n-1})$

sequential: $W(n) = D(n) = O(n)$

parallel: $x = (x_0, \dots, x_{n-1})$

$scan(x)$

if $length(x) = 1$ return (0)

$sums = (x_0 + x_1, \dots, x_{n-2} + x_{n-1})$

$evens = scan(sums) \quad // \quad (0, x_1 + x_2, x_1 + x_2 + x_3 + x_4, \dots)$

$odds = (evens[i] + x[2i] \mid i = 0 \dots \lceil \frac{n}{2} \rceil - 1)$

return $interleave(evens, odds)$

$$W(n) = W(\frac{n}{2}) + O(n) = O(n)$$

$$D(n) = D(\frac{n}{2}) + O(1) = O(\log n)$$

average parallelism: $O(\frac{n}{\log n})$

Reasoning with W and D

Given a DAG with $W(n)$ nodes and depth $D(n)$.

Sequential runtime: $T_1(n) = W(n)$

Time on ∞ processors: $T_\infty(n) = D(n)$

" " " " : $T_p(n) = ?$

$$T_p(n) \geq D(n), \geq \frac{W(n)}{p}$$

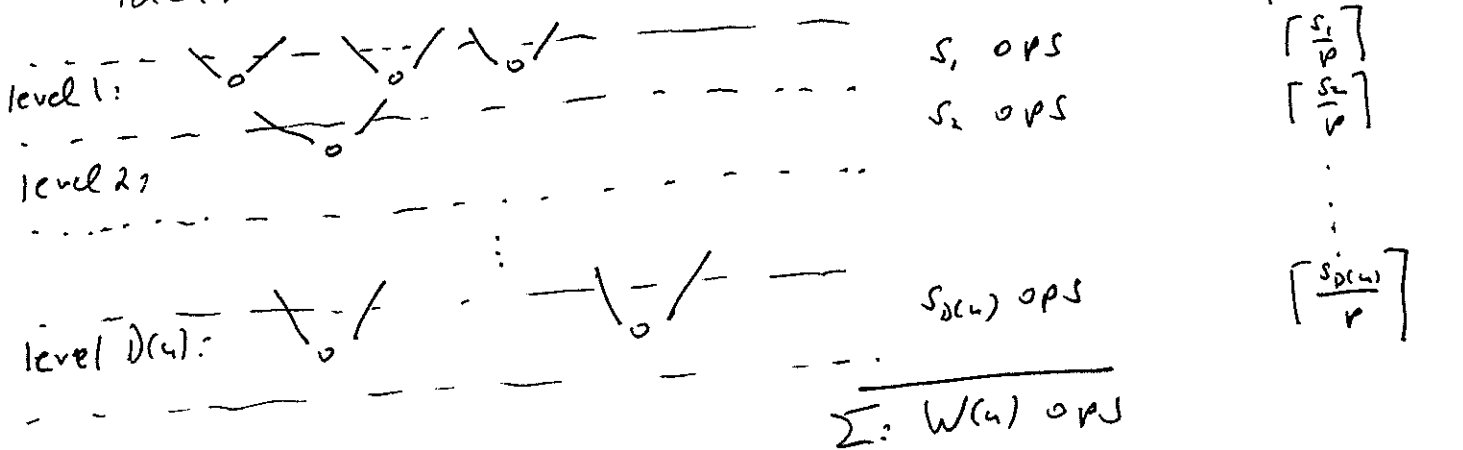
$$T_p(n) \leq D(n) + \frac{W(n) - D(n)}{p}$$

Brent's theorem (1974) $\rightarrow T_p(n) \leq D(n) + \frac{W(n) - D(n)}{p}$ (approximated)

yields: $W(n)/p \leq T_p(n) \leq \frac{W(n)}{p} + D(n)$

Proof of Brent's theorem:

Idea: divide DAG into $D(n)$ levels:



$$T_p(n) \leq \sum_{i=1}^{D(n)} \lceil \frac{s_i}{p} \rceil \leq \sum_{i=1}^{D(n)} \frac{s_i + p - 1}{p} = \frac{1}{p} W(n) + D(n) \frac{p-1}{p}$$

$$= D(n) + \frac{W(n) - D(n)}{p}$$

↑
makes bound non-sharp

Note: level shifts (e.g. node from level $i \rightarrow$ level $i+1$) are not considered

Speedups: $S_p(n) = \frac{W(n)}{D_p(n)}$

$$S_p(n) \leq \frac{W(n)}{D(n)}, \leq P$$

$$S_p(n) \geq \frac{P}{\frac{D(n)}{W(n)} p + 1} \xrightarrow{p \rightarrow \infty} \frac{W(n)}{D(n)}$$

$$S_{\infty}(n) = \frac{W(n)}{D(n)}$$

if n is fixed, speedup is limited. For $n \rightarrow \infty$, speedup can be unbounded.

Strong scaling: Behaviour of $S_p(n)$ for fixed n and increasing p .

Weak scaling: Behaviour of $S_p(n)/p$ for increasing p and increasing n (difficultly proportional to p)

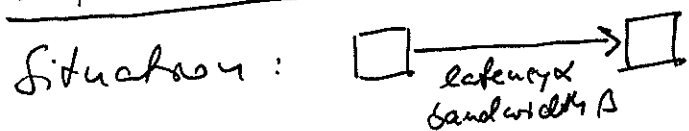
Example: LINPACK benchmark for TOP500 list
→ weak scaling

Example: tree reduction $W(n) = n, D(n) = \lceil \log_2(n) \rceil$

$$S_p(n) \geq \frac{P}{\frac{\lceil \log_2(n) \rceil}{n} p + 1} \xrightarrow{p \rightarrow \infty} \frac{n}{\lceil \log_2(n) \rceil}$$

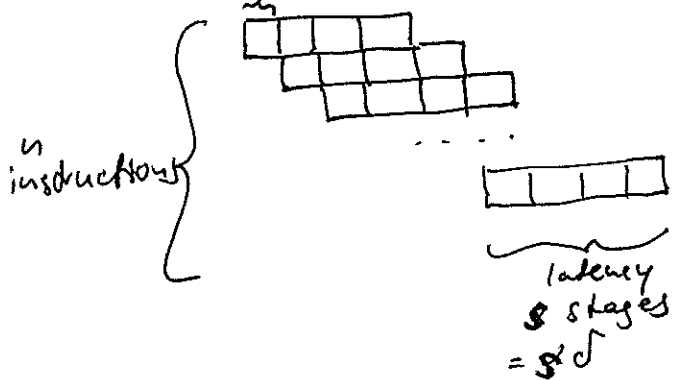
Next: reasoning about latency and through put

α - β Model



How long does it take to send a message of size n ?

Intuition: pipelining of instructions



δ = time per stage
 S = # stages per instruction
 $\beta = 1/\delta$ = bandwidth

Time for n instructions:

$$\begin{aligned} T(n) &= (S-1)\delta + n \cdot \delta \\ &= (S-1)\delta + \frac{n}{\beta} \\ &\approx \text{latency} + \frac{n}{\text{bandwidth}} \end{aligned}$$

Model: $T(n) = \alpha + \frac{n}{\beta}$