

DRIVING FORCE?

Bill McLoch 1
(MSP Note Module Dec. 1992)

Software

Model

Hardware

Architecture-dependent software

Message passing channels
(cf occam)

Distributed Memory Architectures
1980's

Functional programming
(Implicit parallelism)
15 years of research

→ graph reduction
or dataflow

→ Scalable parallel architectures??

Sequential Computing is based on a model

SOFTWARE

VN MODEL

HARDWARE

High level languages
Architecture-independent software

{ Analysis of algorithms
complexity Theory
Performance prediction

{ Focus for technological innovation
Various efficient processor designs

dangers of too much
division re correct model

Parallel computing - mainstream? parallel software industry?

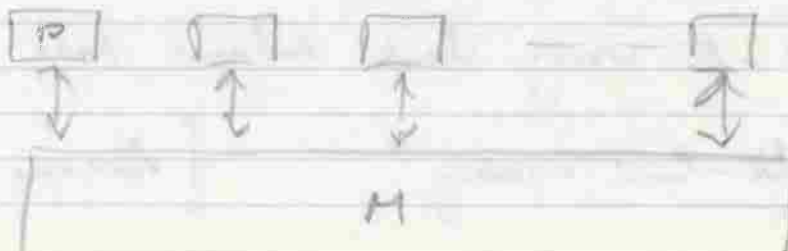
Goals for general-purpose computing

Scalable parallel performance + architecture software independent

1. Idealised parallel computing
2. Special-purpose parallel computing
3. General-purpose

1. Idealised

Synchronisation
current read



Issues for concurrent write

concurrent, nondeterministic, priority value? one succeeds chosen one succeeds

PRAM 1-level memory every memory location is equally far from a processor

Real difficulty not operator, but getting data in right place communication is dominant issue

NC theory NC - $O(\log^k n)$ time on a poly no of processors

Do everything in \mathcal{P} in NC, i.e. do all poly time problems have fast parallel algorithms?

so - complete problems

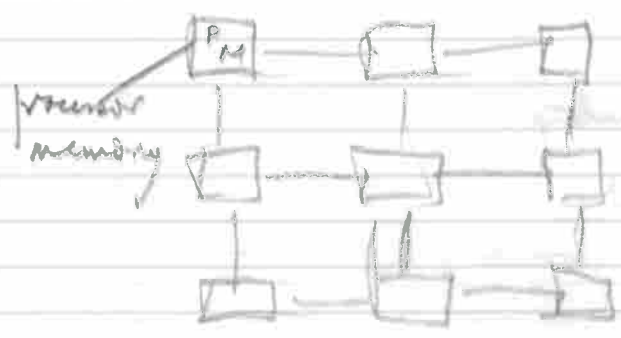
Linear programming
Circuit value problem

General methods for PRAM

Prefix computation, Tree restructuring, List ranking etc

Prove theoretical \rightarrow realistic engineering issues

2. Special-Purpose Parallel Computing



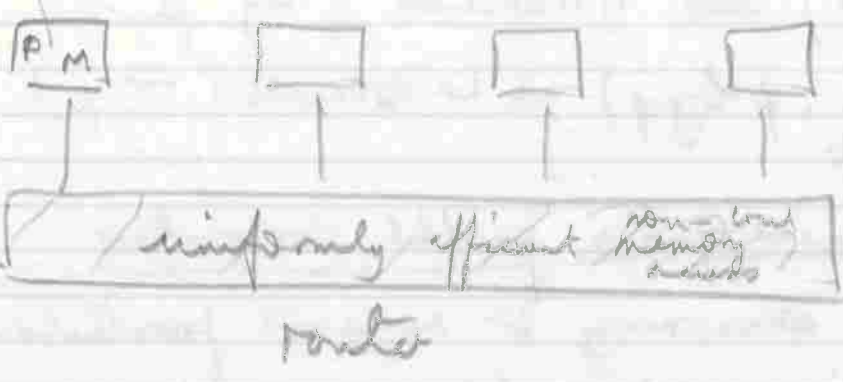
Distributed memory

Algorithm efficiency depends on exploitation of network locality

Have to worry about distance data is transported $\left\{ \begin{array}{l} \text{TRANSMITTER} \\ \text{WIFI chip} \end{array} \right.$

"Multi-level memory" [threads + memory cells require

V. efficient



Virtual shared memory
Single shared address space

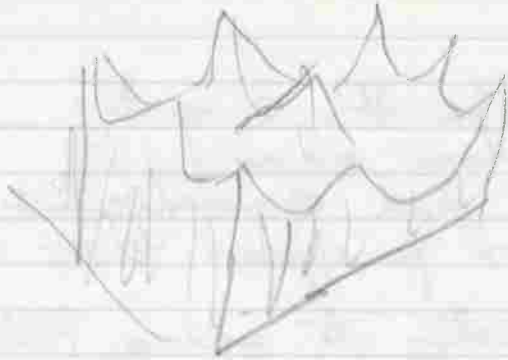
2-level memory

Box in a network e.g. torus wrapped mesh.

Routing box must be able to resolve any formulation efficiently

Randomized routing (Theory end of 70s + Practice)

- Phase 1 send to intermediate node to random node by shortest path
- Phase 2 forward to destination by shortest path



randomized routing
looks out the
to optional points here.

Let h -relation denote problem where each processor sends h packets & receives h .

1-relation can be solved on n -node hypercube (CC etc) in $O(\log n)$ steps (Valiant 1980)



PRAM algorithms with communication overheads can be implemented on hypercube with no overhead.
log n -relation --- (upfal 1984)

PRAM(T, p) \Rightarrow HYPERCUBE ($\log p, p$)

VAUGHAN

PRAM($\frac{T}{\log p}, p \log p$) \Rightarrow HYPERCUBE (T, p)

WPA02

Key ideas for general purpose parallel computing

Fine grain concurrency & maximal parallelism

Parallel structure, latency tolerance, multithreading
(fast context switching)

Fast high-capacity routing networks, communication structures

Uniform memory access.

Problem: retaining locality not a good thing to do!

↳ Look the address space to distribute memory references uniformly

↳ Get around number of threads accessing same pt.

↳ good in theory & practice (PRMA)

BSP model

bulk-synchronous parallel computing

- set of processor memory pairs
- point to point communication network
- mechanism for efficient barrier synchronization

Issue: is it cost-effective to address PRAM implementation

at present time
McLaurin + Valiant think no - hence BSP.

PRAM $\left(\begin{array}{l} g = 1 \text{ no of iterations cycles } \& \text{ radix perm.} \\ g = 2000 \text{ } \\ g = 80, 90? \end{array} \right.$

Challenges

Architecture

Processor design to support a large number of threads
fast context switching, message handling, address

Improved networks

Grand challenge reduce latency lots of ops for every non-local memory access
New optics
"granularity" every op. atomic.

Switch to buses

IO.

Algorithms

Algorithms & complexity theory for BSP

Effect of loading

Sparse, irregular matrices & graphs.

Candidates for GPRC

PRAM, SIMD, dataflow, message passing,
hypercube, graph reduction,

"Maximum range of possibilities occur at about
robustness - endless variations on a few driving models"

BURGESS SHALE

S.J. Gould.

parallel in computation ↔ parallelism in machines
goal of auto-parallelization (1 processor opposite extreme (magic))

consider good get more than you want in parallel computation "parallel slackness"

