Umut Acar
umut.acar@gmail
umut@cmu.edu

Parallelism is <u>HARD!</u>     ?
            is <u>EASY!</u>

Thinking in parallel
design a good parallel arg.
efficiency  < theory
              practice
correctness

learning parallel about as difficult as
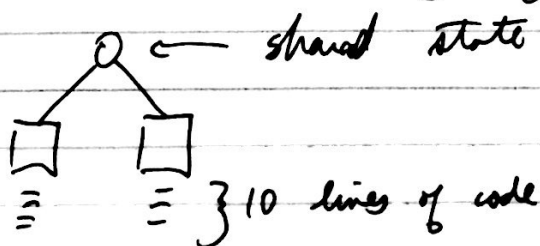sequential algorithms

language-based cost models
    gives
    reason abt complex math. object
    w/o deceiving us

    reflect accurate cost of ~~things~~ reasonably

Key assumption:  pure functional programs

shared state, multiple actors R/W state
    (semantics difficult)
        exponential interleavings of code



        ←— shared state

            ⌒ ⌒ 10 lines of code

how many interleaving?    $2^{20}$    10+10 layout of code
                                       <u>A</u><u>B</u><u>A</u><u>A</u> ....

w/ pure functional (interleaving complexity goes away)

important that it is an abstraction
 · runtime system can use it
 · programmer can't

$\lambda$ calculus - simple cost semantics
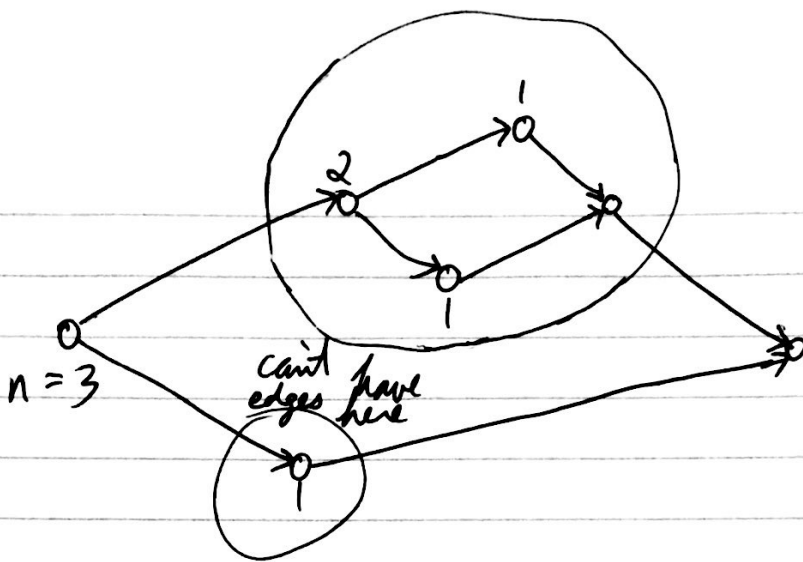 little twist from what you're used to

1. inherently parallel
 $(e, e_2)$ - can eval in parallel
 also write $(e, \| e_2)$

2. twist : cost ?
seq - only care about work = seq. run time
 (additive)
 introduces : S<u>PAN</u> ~> max (not add)

```
fn f x =
  if x ≤ 1 then
       x
  else
     let (a, b) = (f(x-1) ‖ f(x-2))
     in
         a + b
  end
```
 — comma but parallel

$$W(n) = \begin{cases} 1 & \text{if } x \le 1 \\ W(n-1) + W(n-2) + 1 & \text{otherwise} \end{cases}$$

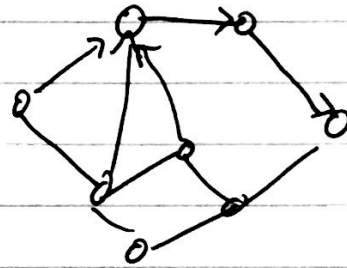$$S(n) = \begin{cases} 1 & \text{if } n \le 1 \\ \max(S(n-1), S(n-2)) + 1 & \text{otherwise} \end{cases}$$
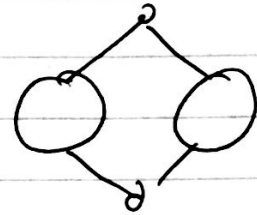
$n = 3$

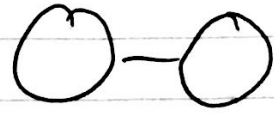can't have edges here

span — longest path



series-parallel DAGs

—or—

nested parallelism.
fork/join parallelism

not DAG for
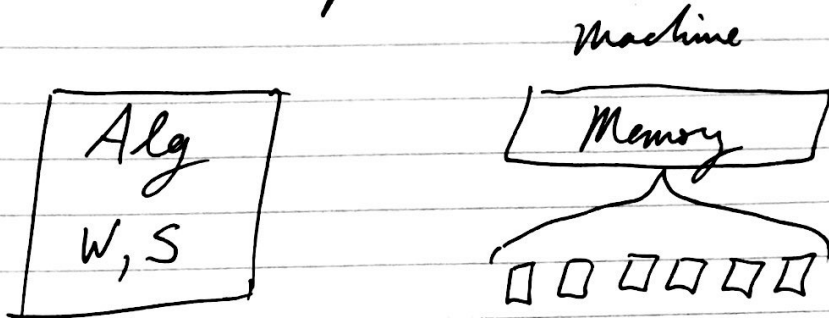in fork/join

(you can w/ futures)

Good parallel alg:

- work close to seq alg.
  (work-efficient)

- low span eg. ($\log^2 n$)

How do we make it truthful
push into runtime system
w/ provably efficient implementation

Provably efficient implementation

Bounded implementation

machine

Alg
W, S

Memory

Want: — result should be the same as sequential

(parallel algs accept seq. semantics)

— take adv of all cores as much
as possible

lower bounds
$$T_P \geq \frac{W}{P}$$
$$T_P \geq S$$

upper bound   caveat: $T_P$ OPT (optimal execution schedule)
is NP-complete!

2 cores

| ~~time~~ | 1 | 2 |
|---|---|---|
| 1 | a | |
| 2 | b | c |
| 3 | d | e |
| 4 | f | |
| 5 | g | |

⎵ schedule

↖

level-by-level scheduler

level (longest distance from source)

not breadth-first (could go to
'g' before 'f' is done
(often 'c')
)

scheduling difficult
make sure all deps are completed
before proceeding on edge in DAG

$W_i$ is total # vertices at level $i$

$$W = \sum_{i=1}^{S} W_i \qquad (S - \text{span})$$

$$T_P = \sum_{i=1}^{S} \left\lceil \frac{W_i}{P} \right\rceil = \sum_{i=1}^{S} \left\lfloor \frac{W_i}{P} \right\rfloor + 1$$

↖ time each level takes

$$= \left( \sum_{i=1}^{S} \left\lfloor \frac{W_i}{P} \right\rfloor \right) + S \;\leq\; \sum_{i=1}^{S} \frac{w_i}{P} + S = \boxed{\frac{W}{P} + S}$$

Brent's Theorem 1972

4

$$T_P \leq \frac{w}{P} + S \leq 2 * OPT$$

$\uparrow$

Why?
$$T_P \geq \frac{w}{P}$$
$$T_P \geq S$$

$$T_P \geq \max\left(\frac{w}{P}, S\right)$$

$$OPT \geq \max\left(\frac{w}{P}, S\right)$$

Justifies Truthfulness of cost model

(W/in factor of two)

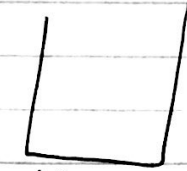The cost of scheduler

vs

The length of schedule ($T_P$)

wasteful
(better!)
(greedy)

$10 \cdot P + 1$

one level only uses one processor

variant: greedy

if one exists idle processor

and exists ready vertex

then assign

guarantees some bound:

$$T_P \leq \frac{w}{P} + S \frac{(P-1)}{P}$$

Pf: arora Blumpe Plato
1997/8



wrk
bucket

idle
bucket

toss coin into work bucket
idle toss coin into idle bucket

each step collect $p$ coins

$$T_p = \frac{\text{total coins}}{p}$$

( why?
collect $p$
at a time )

at end: work bucket has # coins = work

idle bucket:
at any step: $\boxed{\text{max is } p-1}$
(computation not done yet)

for/any step w/ idle processor:

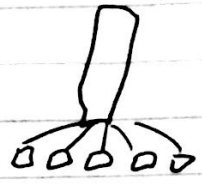span of remaining DAG?
↓
it reduces by $1$
(making progress on span)
must have finished a layer of DAG

total idle coins $(p-1) \times S \implies T_p = \frac{w + (p-1)s}{p}$
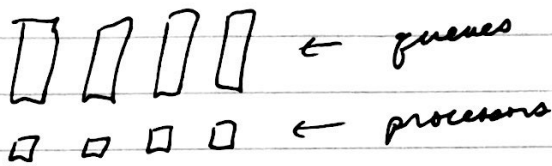
Current exposition does not acct for _cost_
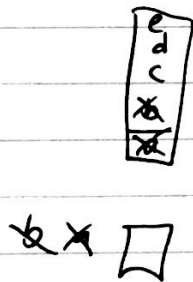of <u>scheduler</u>

Natural way to implement:

— Queue of work ( <u>centralized</u> )
    can only serve work one-at-a-time

— Distributed Queues.

    ← queues

    ← processors

vertex executes — spawns 2 vertices

← manage as a stack

if creates new vertices,
pushes onto its stack

load balancing ⇒ moving vertices to other
                                      processors
processor tries to steal work from
    another processor queue

stealing operates at top end
          they don't contend
               (.

— concurrent algorithm (mutating, effectful, shared state

difficult to implement efficiently, correctly

1995 - 2005
↑ reasonable

w/o language model, you effectively implement your own scheduler

this alg gives $E[T_P] = \frac{W}{P} + S$ in expectation

↑ steals happening uniformly

actually includes cost of scheduler

high-level cost model can be truthful

implemented many times (greedy scheduler)

work-stealing

Berton - Sleep

Leisserson
Blelloch
Blumofe

has good cache locality

simulates seq. execution except at steal points

Java, Parallel ML, Silk, Haskell