# CPSC 465/565 Theory of Distributed Systems

James Aspnes

2023-09-01

# Today's exciting topics

- Message passing complexity measures
- Broadcast and convergecast

# Message passing model



- Buffer $b_ij$ holds messages from $p_i$ to $p_j$
- Execution: $C_0\alpha_1 C_1\alpha_2 C_2 \ldots$
    - $\alpha_t =$ delivery event at step $t$
- Nondeterminism: adversary chooses next event from enabled events
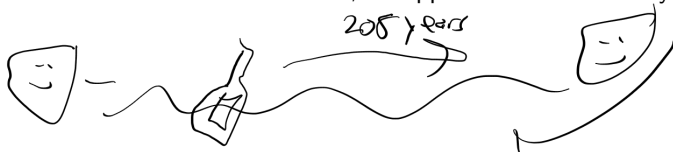- Fairness: Every message sent is eventually received

What about complexity?

# Complexity measures for message-passing

Last time:

- Local computation: free!
  - Justification: usually fast relative to message delivery times
- Message complexity = number of messages sent
- Bit complexity = total size of all messages

This time:

- Time complexity = time to finish
  - But: even in fair executions, no upper bound on delivery time
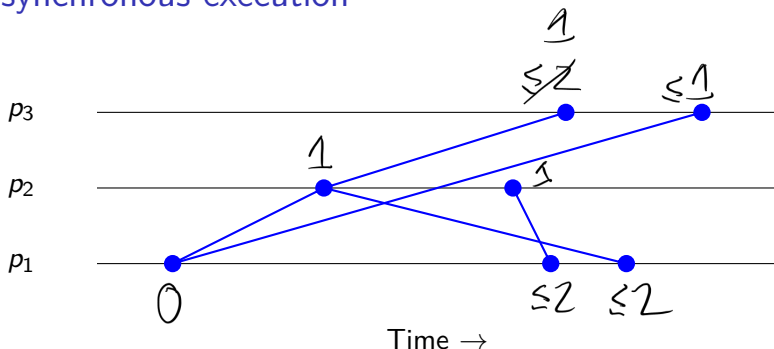
# Time complexity without time bounds



- *Define* maximum message delay = 1 time unit.
- *Define* gap between two local computation events = 0 time units.
- Assign each event the max time consistent with these assumptions.

This gives a time measure that is agnostic about actual time delays.

(Similar to assumption behind $O()$ notation.)

# An asynchronous execution



Time $\rightarrow$

Assign highest possible time to each event consistent with:

► $T(\text{first step}) = 0$.
► $T(\text{receive}) \leq T((\text{send})) + 1$.
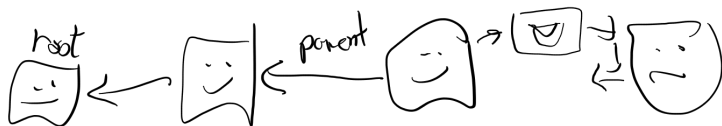► Consecutive local computation events have same time.

# Example: flooding with parent pointers

```
1  initially do
2  |   if pid = root then
3  |   |   parent ← root
4  |   |   send m to all neighbors
5  |   else
6  |   |   parent ← ⊥

7  upon receiving m from j do
8  |   if parent = ⊥ then
9  |   |   parent ← p
10 |   |
11 |   |   send m to all neighbors
```

Builds a spanning tree with broadcast as side-effect.

# Correctness



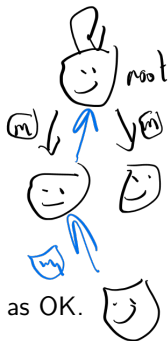Claim: eventually parent pointers form rooted spanning tree.

- Safety:
    1. $parent_i = \bot$ or $parent_i$ is node on path to root
    2. If $m \in b_{ij}$, then $parent_i \neq \bot$.
- Liveness: Eventually $parent_i \neq \bot$.

# Safety proof



1. $\text{parent}_i = \bot$ or $\text{parent}_i$ is node on path to root
2. If $m \in b_{ij}$, then $\text{parent}_i \neq \bot$.

▶ Check initial configuration:
   1. $\text{parent}_i = \bot$ unless $i = \text{root}$; count $\text{parent}_{\text{root}} = \text{root}$ as OK.
   2. If $m \in b_{ij}$, then $i = \text{root}$ and $\text{parent}_{\text{root}} \neq \bot$.
▶ Check each claim against steps that could make it false:
   1. If $i$ sets $\text{parent}_i$ to $j$:
      ▶ $i$ receives $m$ from $j$
      ▶ $m \in b_{ij} \Rightarrow \text{parent}_j \neq \bot \Rightarrow \text{parent}_j$ is on path to root
   2. If $i$ sends $m$:
      ▶ Same code sets $\text{parent}_i \neq \bot$
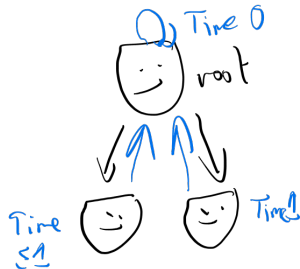
# Liveness proof

Eventually $parent_i \neq \perp$.

Same induction on distance as simple flooding.

But now we will include time.

- ▶ Claim: For all $i$, $parent_i \neq \perp$ by time $d(i, \text{root})$.
- ▶ Proof: By induction on $d(i, \text{root})$:
  - ▶ Base: $d(i, \text{root}) = 0 \Rightarrow i = \text{root}$, so $parent_i = i$ at time 0.
  - ▶ Induction step: Suppose $d(i, \text{root}) = d + 1$
    - ▶ Then there is some neighbor $j$ that sets $parent_j \neq \perp$ by time $d$.
    - ▶ When $j$ does this, it sends $m$ to $i$.
    - ▶ At most one time unit later ($\leq d + 1$), $i$ receives $m$.
    - ▶ When $i$ receives $m$, sets $parent_i \neq \perp$ if not already set.

# Time complexity

$D = $ diameter
$= $ max dist
between 2 nodes

Count until last message delivered.

- Most distant node is at $\max_i d(i, \text{root}) \leq \text{diam}(G)$.
- Every node sets $\text{parent}_i \neq \bot$ by time $D$.
- Last message is sent by time $D$.
- Last message is received by time $D + 1$.

So time complexity is $D + 1$.

## Detecting termination

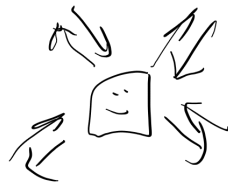Add acknowledgments:

```
1  upon receiving m from j do
2  │   if parent = ⊥ then
3  │   │   parent ← j
4  │   │   send m to all neighbors
5  │   else
6  │   │   send nack to j

7  upon receiving ack from j do
8  │   children ← children ∪ {j}

9  upon receiving nack from j do
10 │   nonChildren ← nonChildren ∪ {j}

11 as soon as children ∪ nonChildren = neighbors do
12 │   send ack to parent
```



Done when root gets ack or nack from all neighbors.

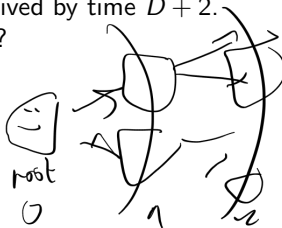# Complexity of the terminating version

- Message complexity
    - Every edge gets $m$ in each direction
    - Plus ack or nack in response
    - Total is $4|E|$ messages.
- Time complexity
    - Outgoing messages reach everybody in time $D$.
    - All nacks are sent by time $D + 1$.
    - All nacks are received by time $D + 2$.
    - What about acks?

$2|E|$
recruiting
msgs

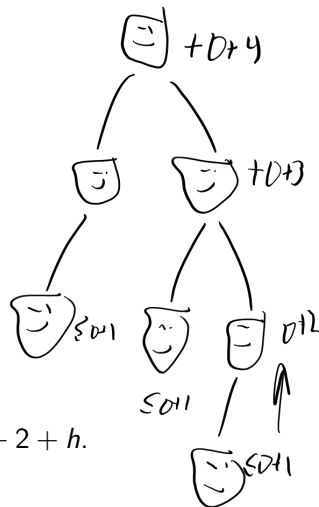$+2|E|$
ack $\&$ nacks

$D$

# Timing of acks

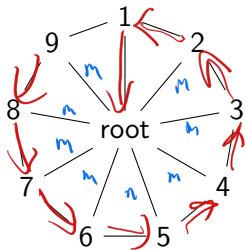Look at position of each node in the tree:

- ▶ If I am a leaf:
  - ▶ I send ack once I collect all nacks
  - ▶ This happens by time $D + 2$
- ▶ If I am a parent:
  - ▶ I have to wait for all my kids.

Induction hypothesis:

- ▶ Node at height $h$ sends ack at time $\leq D + 2 + h$.
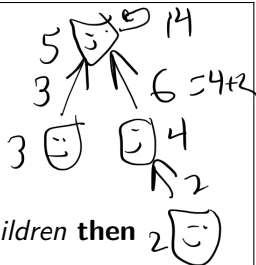- ▶ $\Rightarrow$ finish by time $D + 2 +$ (height of tree).

# How tall can the tree get?



- Max distance from root is 1, so all nodes recruited by time 1.
- But maybe some nodes are recruited faster.
- Worst case: height of tree = $n - 1$.
- Gives $D + 2 + (n - 1)$ worst-case time for broadcast with acks.

## Convergecast

Sometimes it's convenient to collect data from an existing tree:



**1** **initially do**
**2**     **if** *I am a leaf* **then**
**3**        send input to parent

**4** **upon receiving** $v_c$ **from** $c$ **do**
**5**     append $(c, v_c)$ to buffer
**6**     **if** buffer *contains messages from all my children* **then**
**7**        $v \leftarrow f(\text{buffer}, \text{input})$
**8**        **if** pid $=$ root **then**
**9**           **return** $v$
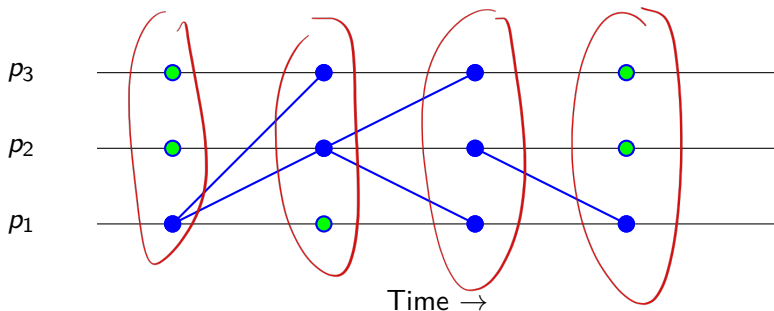**10**        **else**
**11**           send $v$ to parent

**Broadcast** to build the tree and **convergecast** to collect data.

# Performance of convergecast

- Message complexity $= n - 1$
- Time complexity:
    - Leaves send at time 0
    - Parent sends by max child time $+1$
    - $\Rightarrow$ each $i$ sends by height($i$)
    - $O(n)$ time in worst case.

# Synchronous model

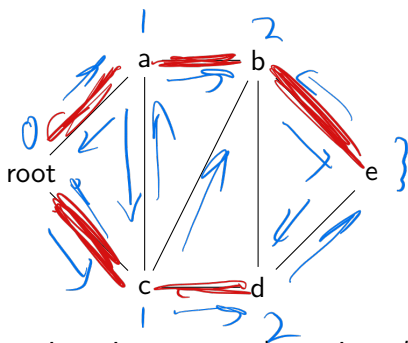

Assumption: every message delay $= 1$ exactly.

Formal version:

- Sequence of rounds $r = 0, 1, 2, \ldots$
- Every node gets one delivery event per round.
- Every message sent in $r$ is received in $r + 1$.

# Synchronous broadcast gives shortest-path tree



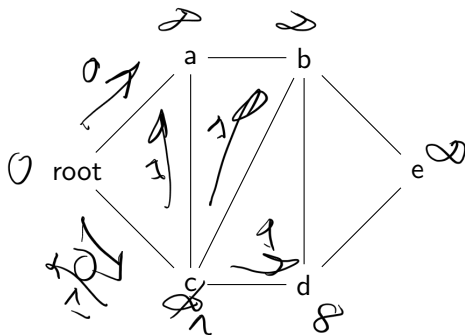Proof: By induction, $i$ receives $m$ exactly at time $d(\text{root}, i)$.

end center

# What if we don't have synchrony?

Attempt 1: Use Dijkstra's algorithm.

```
1  initially do
2  |   if pid = initiator then
3  |   |   dist ← 0
4  |   |   send dist to all neighbors
5  |   else
6  |   |   dist ← ∞
7  upon receiving d from p do
8  |   if d + 1 < dist then
9  |   |   dist ← d + 1
10 |   |   parent ← p
11 |   |   send dist to all neighbors
```

# Dijsktra in action



Claim: Every process converges to correct distance.

▶ Safety: $\text{dist}_i \geq d(i, \text{root})$ and for any $d \in b_{ij}$, $d \geq d(i, \text{root})$.
▶ Liveness: $\text{dist}_i = d(i, \text{root})$ no later than time $d(i, \text{root})$.
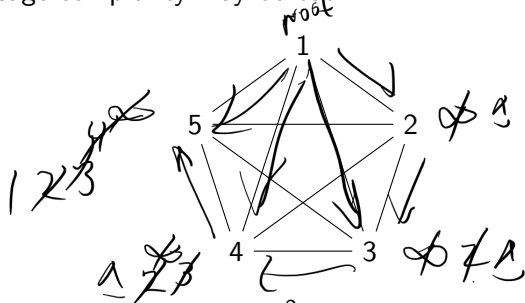
Ind on $d$: base: true for $d=0$
Ind step: closer nbr at dist $d$ sets by $t=d$
$\Rightarrow$ msg $\Rightarrow$ I col to $d+1$ by time $d+1$.

# Worst-case complexity

Liveness proof gives $O(D)$ time.

But message complexity may be bad:



- Careful schedule yields $\Theta(n^2)$ updates.
- $\Rightarrow \Theta(n^3)$ messages.
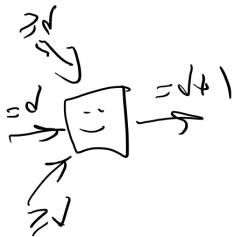- With synchrony, only $2|E| = O(n^2)$ messages.

# Fix by simulating synchrony ($\propto$ synchronizer)

```
 1  initially do
 2  │  if pid = root then
 3  │  │    dist ← 0
 4  │  │    send [= 0] to all neighbors
 5  │  else
 6  │  │    send [≥ 0] to all neighbors

 7  as soon as I received [≥ d] from all neighbors do
 8  │  send [≥ d + 1] to all neighbors

 9  as soon as I received [= d] from at least one neighbor and
    [≥ d] from the others do
10  │  dist ← d + 1
11  │  send [= d + 1] to all neighbors
12  │  halt
```

$(= d] \rightarrow$ I am exactly $d$

$(\geq d] \rightarrow$ I am at least $d$

# Correctness



- Safety:
  ∪
  - If $i$ sends $[\geq d]$, $d(i, \text{root}) \geq d$.
  - If $i$ sends $[= d]$, $d(i, \text{root}) = d$.

  + similar constraints on msgs.

- Liveness:
  - If $d(i, \text{root}) = d$, then
    - For each $d' < d$, $i$ sends $[\geq d']$ by time $d'$.
    - $i$ sends $[= d]$ by time $d$.
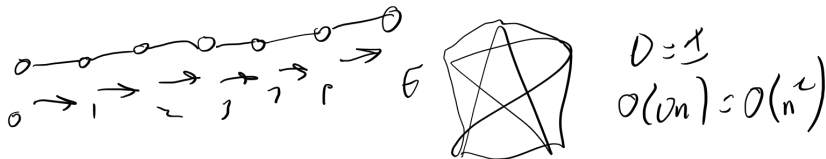
    Easy induction on $d$.

    (rhymes w/ lazy)

# Cost

- From safety property:
    - Distance-$d$ node sends only $[\geq 0], [\geq 1], \ldots, [\geq d-1], [= d]$.
    - This gives at most $D + 1$ messages per direction per edge.
    - $\Rightarrow$ Message complexity $= O(D|E|)$.
- Liveness property:
    - Most distant node sends $[= d]$ by time $d$.
    - $\Rightarrow$ Time complexity $= O(D)$.

For example, get $O(n^2)$ messages and $O(1)$ time on $K_n$.

# Next time

Leader election!