# EC504 ALGORITHMS AND DATA STRUCTURES FALL 2020 MONDAY & WEDNESDAY 2:30 PM - 4:15 PM

Prof: David Castañón, dac@bu.edu

GTF: Mert Toslali, toslali@bu.edu

Haoyang Wang: haoyangw@bu.edu

Christopher Liao: cliao25@bu.edu

**Hw 3. Next week** ←

**· Sw hw  2 tomorrow** ←

**→ Exam 1 early march** ←

———————→ | Graphs →

Exam 1: Hws 1-3.
Format: ?
    Yes.   Lecture time.
    No.
Similar to Hw + Quizzes
No coding.
Open book & Notes.
Not open Google . . .

# Hash Tables

o Hashing
  o Technique supporting insertion, deletion and search in average-case constant time
  o Operations requiring elements to be sorted (e.g., FindMin) are not efficiently supported

o Generalizes an ordinary array,

  o Key property: direct addressing

  o An array is a direct-address table: Key value is position of data in array

o Main idea: Transform key into index, compute the index, then use an array of size N
  o Key k: data stored at h(k)  (hashing)

o Basic operation is  in O(1)!

# Hash Function

o Mapping from key to array index is called a hash function

    o Typically, many-to-one mapping
    o Different keys map to different indices
    o Distributes keys evenly over table

o Collision occurs when hash function maps two keys to the same array index

# Collision Resolution - 2

o  Open addressing
  o  If slot is busy, design sequence of other slots to be searched
  o  probe alternative cell $h_1(K), h_2(K), \ldots,$ until an empty cell is found.
  o  $h_i(K)$ = (hash(K) + f(i)) mod m, with f(0) = 0
  o  f: collision resolution strategy

o  Several approaches
  o  Linear Probing:  f(k) = k  ✔
  o  Quadratic Probing: f(k) = $k^2$  ✔
  o  Double Hashing: two hash functions
  o  Cuckoo Hashing: (more to come)

# Linear Probing

o f(i) =i
  o cells are probed sequentially (with wrap-around)
  o $h_i$(K) = (hash(K) +  i) mod m

o Insertion:
  o Let K be the new key to be inserted, compute hash(K)
  o For i = 0 to m-1
    o compute L = ( hash(K) + I ) mod m
    o If T[L] is empty, then we put K there and stop
  o If we cannot find an empty entry to put K, it means that the table is full and we should report an error: Table is full

o Problem:  We no longer have O(1) find, insert for worst case.

# Example

o E.g, inserting keys 89, 18, 49, 58, 69 with hash(K)=K mod 10 (not prime!)

   o $h_i$(K) = (hash(K) + i) mod m

o

| | Empty Table | After 89 | After 18 | After 49 | After 58 | After 69 |
|---|---|---|---|---|---|---|
| 0 | | | | 49 | 49 | |
| 1 | | | | | 58 | 58 |
| 2 | | | | | | 69 |
| 3 | | | | | | |
| 4 | | | | | | |
| 5 | | | | | | |
| 6 | | | | | | |
| 7 | | | | | | |
| 8 | | | 18 | 18 | 18 | 18 |
| 9 | | 89 | 89 | 89 | 89 | 89 |

# Quadratic Probing

o Two keys with different home positions will have different probe sequences
- o e.g. m=101, h(k1)=30, h(k2)=29
- o probe sequence for k1: 30,30+1, 30+4, 30+9
- o probe sequence for k2: 29, 29+1, 29+4, 29+9

o If the table size is prime, then a new key can always be inserted if the table is at least half empty (see proof in text book)
- o Secondary clustering
  - o Keys that hash to the same home position will probe the same alternative cells
  - o Simulation results suggest that it generally causes less than an extra half probe per search
  - o To avoid secondary clustering, the probe sequence need to be a function of the original key value, not the home position

# Quadratic probing

o E.g, inserting keys 89, 18, 49, 58, 69 with hash(K)=K mod 10 (not prime!)

    o $h_i(K) = (hash(K) + i^2)$ mod m

o

| | Empty Table | After 89 | After 18 | After 49 | After 58 | After 69 |
|---|---|---|---|---|---|---|
| 0 | | | | 49 | 49 | 49 |
| 1 | | | | | | |
| 2 | | | | | 58 | 58 |
| 3 | | | | | | 69 |
| 4 | | | | | | |
| 5 | | | | | | |
| 6 | | | | | | |
| 7 | | | | | | |
| 8 | | | 18 | 18 | 18 | 18 |
| 9 | | 89 | 89 | 89 | 89 | 89 |

204

Quadratic Probing

# Double Hashing

o To alleviate the problem of clustering, the sequence of probes for a key should be independent of its primary position => use two hash functions: hash() and hash2()

   o f(i) = i * $hash_2(K)$

   o E.g. $hash_2(K)$ = R - (K mod R), with R is a prime smaller than m


o $hash_2(K)$ must never evaluate to zero

   o For any key K, $hash_2(K)$ must be relatively prime to the table size m. Otherwise, we will only be able to examine a fraction of the table entries.

   o One solution is to make m prime, and choose R to be a prime smaller than m

205

# Double Hashing

o E.g, inserting keys 89, 18, 49, 58, 69 with hash(K)=K mod 10 (not prime!)

   o $h_2$(K) = (hash(K) + i * $h_2(K)$) mod m, where $h_2(K) = K$ mod 7

o

| | Empty Table | After 89 | After 18 | After 49 | After 58 | After 69 |
|---|---|---|---|---|---|---|
| 0 | | | | | 58 | 69 |
| 1 | | | | | | |
| 2 | | | | | | |
| 3 | | | | | ~~58~~ | ~~58~~ |
| 4 | | | | | | 69 |
| 5 | | | | | | |
| 6 | | | | 49 | 49 | 49 |
| 7 | | | | | | |
| 8 | | | 18 | 18 | 18 | 18 |
| 9 | | 89 | 89 | 89 | 89 | 89 |

58 mod 7 = 2.

6.
69

206

# Rehashing ✓

o When hash table is near full, increase size, rehash all entries
  o Amortized still O(1)…

o When to rehash
  o When table is half full (λ = 0.5) ✓
  o When an insertion fails
  o When load factor reaches some threshold ✓
  o Works for chaining and open addressing ✓

o

$n/m$ = load factor.

# Big Problem with Open Address: Deletion

o If you delete, you break the chain of insertions and lose ability to find!

    o Solution: Add an extra bit to each table entry, and mark a deleted slot by storing a special value DELETED

o $hash_2(K)$ must never evaluate to zero

    o For any key K, $hash_2(K)$ must be relatively prime to the table size m. Otherwise, we will only be able to examine a fraction of the table entries.

    o One solution is to make m prime, and choose R to be a prime smaller than m

o

# Perfect Hashing

o Choose a hash function with no collisions: Hard!

    o The expected cost of a lookup in a chained hash table is $O(1 + \alpha)$ for any load factor $\alpha$

    o Expected cost of a lookup in these tables is not the same as the expected worst-case cost

o Theorem: Assuming truly random hash functions, the expected worst-case cost of a lookup in a linear probing hash table is $\Omega(\log n)$.

o Theorem: Assuming truly random hash functions, the expected worst-case cost of a lookup in a chained hash table is $\Theta(\log n / \log \log n)$.

    o Proofs: CLRS, exercise 11-1 and 11-2.

209

# Perfect Hashing - 2

o Bottom line: perfect hashing needs O(1) inserts, lookups
  o Chaining and linear probing are a long way from this ✓
  o This is for expected worst case, not actual worst case (which is worse!)
  o Need creative techniques to approach perfect hashing

o Let's try a new idea: Cuckoo Hashing ✓

# Cuckoo Hashing

o Cuckoo hashing is a simple hash table where
  o Lookups are worst-case O(1).
  o Deletions are worst-case O(1).
  o Insertions are amortized, expected O(1).
  o Insertions are amortized O(1) with reasonably high probability.
o Key idea: Maintain two table, each with m elements
  o Two hash functions $h_1(K), h_2(K)$
  o Every element K will be either in position $h_1(K)$ in the first table or $h_2(K)$ in the second table

# Cuckoo Hashing

o Lookups are O(1) because only 2 locations must be searched

o Deletions take O(1) because only 2 locations must be searched

o To insert an element, try placing it in $h_1(K)$. If empty, place it there

o If not empty, place it there, and kick out existing one (J) and try in $h_2(J)$

o If $h_2(J)$ is busy, place J there, kick out L and try placing it in $h_1(L)$

o Repeat, alternating, until items stabilize

| $T_1$ |
|-------|
|       |
| 32    |
|       |
| 84    |
| 59    |
|       |
| 93    |
| 58    |
|       |
|       |

| $T_2$ |
|-------|
| 97    |
|       |
| 26    |
|       |
| 41    |
|       |
| 23    |
|       |
| 53    |
|       |

# Cuckoo Hashing Example



213

# Cuckoo Hashing

o What can go wrong?
  o We run into a cycle!

# Cuckoo Hashing

o What can go wrong?
  o We run into a cycle!

o If that happens, perform a rehash by choosing two new hash functions and inserting all elements back into the tables

o Multiple rehashes might be necessary before this succeeds

o Cycles only arise if we revisit the same slot with the same element to insert

215

# Cuckoo Hashing Results

o Hard probabilistic analysis: based on random bipartite graphs
  o Uses the Cuckoo graph
  o Beyond scope of our course — still a topic for research
o m: size of one of the hash tables.  n:the number of edges between the two tables

o Theorem: If m = (1 + ε)n for ε > 0, the probability that cuckoo graph contains a complex connected component is O(1 / m)

  Keep a little over half the cells empty: $\dfrac{n}{2m} = \dfrac{1}{2 + 2\epsilon}$

o Theorem: The expected, amortized cost of an insertion into a cuckoo hash table is O(1 + ε + ε-1)

*(handwritten annotations)*

Problem: Complex Cycle. → more tha
One single cycle. P[ 1 ~ 1/(m)].

*(handwritten table/diagram annotations)*

| $T_1$ | | $T_2$ |
|---|---|---|
| 19 | 19 / 16 | 16 |
| 91 | 91 | 58 |
| 32 | 8 | |
| 10 | 88 | 75 |
| 6 | 6 | |
| 53 | 53 | 88 |
| 26 | 26 | |
| | 4 | |

216

$$\frac{m}{n} = \frac{50100}{.50000} + 17 \checkmark \qquad \frac{n}{2m} \leqslant 0.5$$

# Cuckoo Hashing Variations

o The hash functions chosen need to have a high degree of independence for results to hold
   o Once numbers of keys gets close to 1/2, failure is imminent!

o Cuckoo hashing with k ≥ 3 tables tends to perform much better than Cuckoo hashing with k = 2 tables
   o With k = 3, you can load tables up to 90% before you run into cycles with enough probability

o Another idea: slots in a cuckoo hash table can store multiple elements
   o When displacing an element, choose a random one to move and move it.
   o Works well, makes it unlikely to have long chains ✓

# Cuckoo Hashing

o Tricky to analyze
  o Everything moves around, two tables, change hash functions, reinsert

o If that happens, perform a rehash by choosing two new hash functions and inserting all elements back into the tables

o Multiple rehashes might be necessary before this succeeds

o Cycles only arise if we revisit the same slot with the same element to insert

# Evolution of Hashing

- Hashing is fast, O(1) inserts, deletes, find
  - But lacks ability for successor, predecessor, find min, ...O(n)!

- New data structures developed for fast operations
  - If data range limited to [0,U]: can do bitwise orderings (similar to Radix sort)
  - Van Emde Boas trees
  - Structures using **tries**: we'll discuss later
  - x-fast tries: uses cuckoo hashing together with bitmap tries: find in O(1), insert in O(log(U)), find min O(log(log(U)))
  - y-fast tries: x-fast trie on top of forest of red-black trees with all operations O(log(log(U)))

219

# Fibonacci Heaps

- Fredman-Tarjan 1986
  - CLRS Chapter 19 ✓

- Binary heaps:  Insert:  O(log(n)), Merge: O(n), DeleteMin: O(log(n));
  DecreaseKey: O(log(n))
- Lazy Binomial heaps: Insert: O(1); Merge: O(1); DeleteMin: O(log(n)) (amortized)
  DecreaseKey: O(log(n)) ✓

- Network optimization algorithms require priority queues where keys are
  decreased much more often than inserts, deletes, pop minima
  - Can we find a data structure to make DecreaseKey O(1) amortized?

Lazy:   A Book of 5 Rings   Mushashi:
"Do nothing which is of no use."

# DecreaseKey in Binomial Heaps

- Assume you can find the in the Binomial Heap in O(1) (may require a dictionary or hashmap)
  - Decrease key and up-heap it  (O(log(n)))
  - If key is at root of subtree, update minimum pointer.



$O(\log n)$

STL...

# DecreaseKey in Binomial Heaps

*Ker*

- Assume you can find the in the Binomial Heap in O(1) (may require a dictionary or hashmap)
  - Decrease key and up-heap it  (O(log(n)))
  - If key is at root of subtree, update minimum pointer.

# An Unusual Idea

- Cut subtrees if heap order violated: O(1)
  - Loses binomial tree structure...
  - But we may be able to get by with that



not binomial

log(a) height

6

# Fibonacci Heaps

- Similar to Binomial Heaps ✓
    - Forest of heap-ordered trees, but trees do not have to be in binomial shape
    - Defer consolidation as in lazy binomial trees ✓
    - DecreaseKey by breaking off the subtree and adding subtree to forest of trees
- Binary heaps:  Insert:  O(log(n)), Merge: O(n), DeleteMin: O(log(n));
  DecreaseKey: O(log(n))

- Lazy Binomial heals: Insert: O(1); Merge: O(1); DeleteMin: O(log(n)) (amortized)
  DecreaseKey: O(log(n))

- Problem: $\Theta(n)$ number of trees in worst case...need to manage complexity

# Fibonacci Heaps

- Notation: Order of a node = number of children
  - In binomial trees with root of order h, there are $2^h$ nodes
  - If cut trees are no longer binomial, they may have fewer keys

- e.g. number of nodes $\Theta(k^2)$, number of trees $\Theta(k)$
  - Need to avoid this! Want number of nodes exponential in number of trees
  - Must impose some structural constraints



8

# Fibonacci Heaps

- Structural constraint: limit number of cuts of children to a non-root node
  - At most one cut can be done without restructuring
  - Mark node that has lost one child
  - If a non-root node loses a second child, we cut the node from its parent also
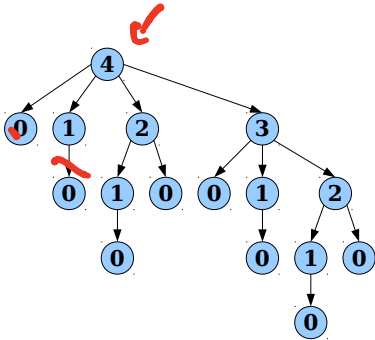    - May be recursive...

# Fibonacci Heaps

- Structural constraint: limit number of cuts of children to a non-root node
  - At most one cut can be done without restructuring
  - Mark node that has lost one child
  - If a non-root node loses a second child, we cut the node from its parent also
    - May be recursive…

# Fibonacci Heaps

- Structural constraint: limit number of cuts of children to a non-root node
  - Mark node that has lost one child
  - If a non-root node loses a second child, we cut the node from its parent also
    - May be recursive...

# Fibonacci Heaps

- Cut operation: cut node v from parent p
  - Unmark v. Cut v from p
  - If p is not marked and is not the root of a tree, mark it
  - If p was already marked, recursively cut p from its parent

- If we do a few decrease-keys, then the tree won't lose "too many" nodes.
  - If we do many decrease-keys, the information slowly propagates to the root

- DeleteMin: complexity O(h), where h is height of tallest tree
  - In Binomial heaps, $h \in O(\log(n))$
  - What is it now?

Can't happen

12

# Fibonacci Heaps

- Minimum number of nodes in tree of rank h



Ranks of
nodes shar.

# Fibonacci Heaps

- Lemma: Number of keys in a tree of order k (root rank k) is $F_{k+2}$, the k+2 Fibonacci number

  - Implies that number of keys grows exponentially with root rank: $\left(\dfrac{1 + \sqrt{5}}{2}\right)^k$

- Max rank O(log(n)), height is no larger than rank



| $T_0$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_5$ |
|-------|-------|-------|-------|-------|-------|
| $F_2 = 1$ | $F_3 = 2$ | $F_4 = 3$ | $F_5 = 5$ | $F_6 = 8$ | $F_7 = 13$ |

# Fibonacci Heaps

- Operations
  - Insert: O(1), just add another tree to heap with insert node, update min
  - Merge: O(1), just link the two sets of trees, update min
  - Find_Min: O(1), just read it
  - Delete_Min: Delete min root, consolidate trees of the same rank. Analysis identical to Binomial Heap, O(log(n)) amortized
  - Decrease_Key: ???

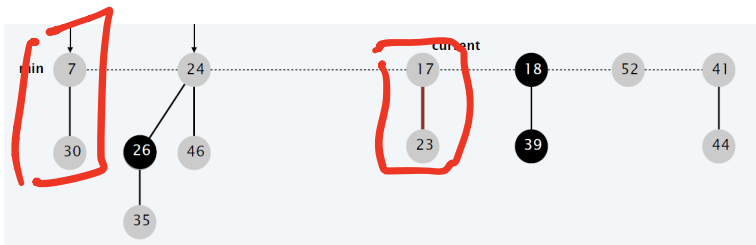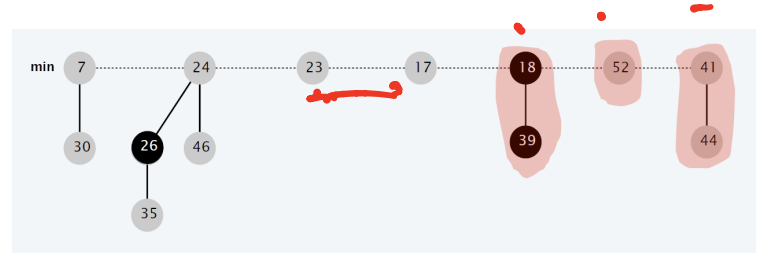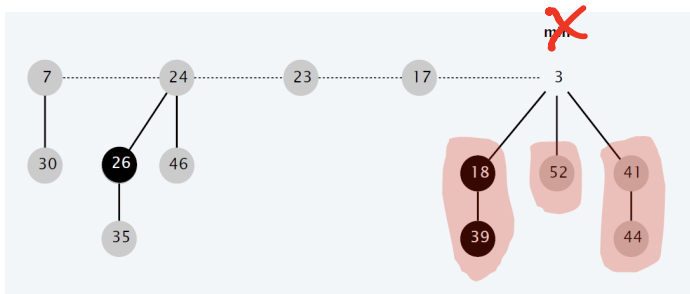- Amortized analysis: Potential function $\Phi(\mathscr{H})$ = number of trees + 2 * number of marked nodes

# Amortized Analysis of DecreaseKey

- Amortized analysis: Potential function $\Phi(\mathcal{H})$ = number of trees + 2 * number of marked nodes



n = 14    rank(H) = 3    trees(H) = 5    marks(H) = 3

5 + 6 = $\Phi$ = 11

# Amortized Analysis of DeleteMin

- Amortized analysis: Potential function $\Phi(\mathcal{H})$ = number of trees + 2 * number of marked nodes
  - Insert: increases potential by 1, O(1) work so amortized O(1)
  - Delete_min:
    - Promoting children of min root increases trees by max rank in $\mathcal{H}$: rank($\mathcal{H}$)
    - When heap $\mathcal{H}$ has k trees, consolidate is $\Theta(k)$+ rank($\mathcal{H}$)
    - Number of trees after consolidation: less than or equal to rank($\mathcal{H}'$) + 1
      - No repeated ranks
    - May lose some marks
    - Amortized cost: O(rank($\mathcal{H}$)) + O(rank($\mathcal{H}'$)), the latter of which is O(log(n))
    - Claim: O(rank($\mathcal{H}$)) is O(log(n)) also!

# Amortized Analysis of DeleteMin

- DeleteMin
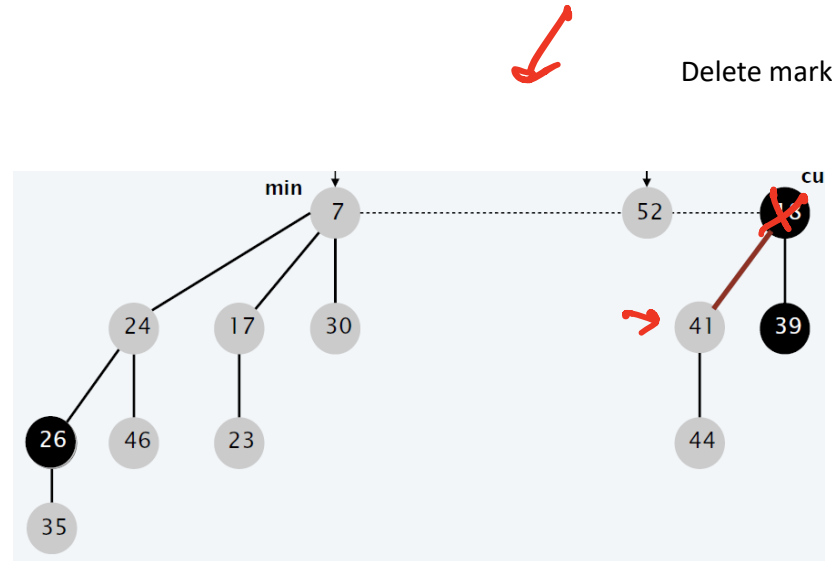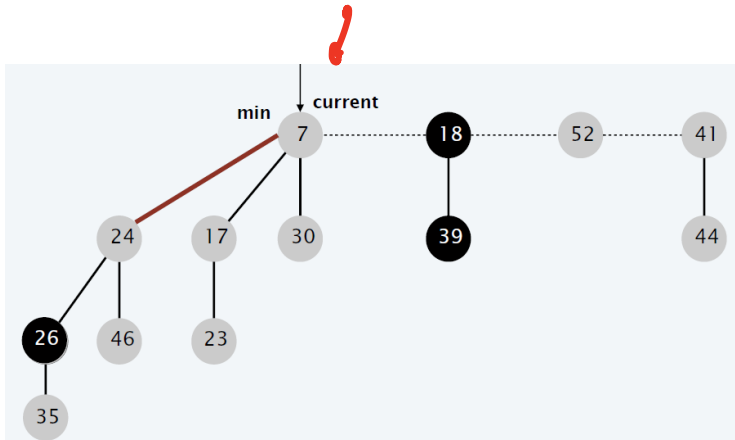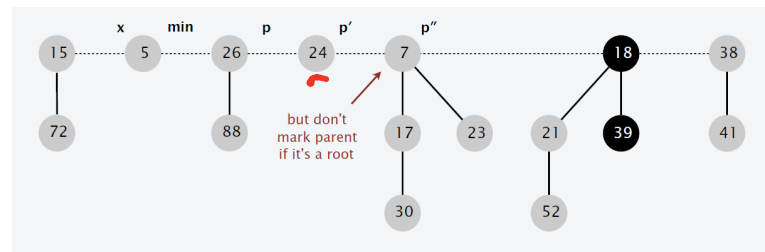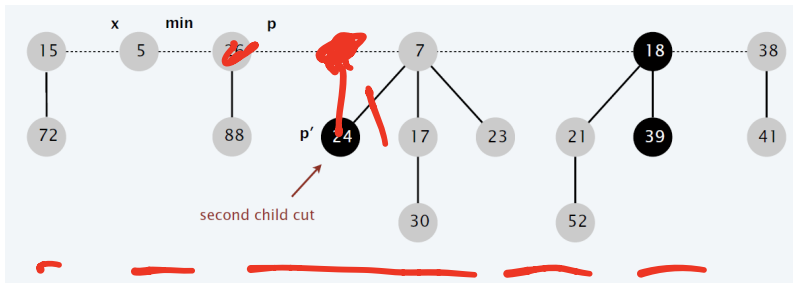
# Amortized Analysis of DeleteMin

- DeleteMin

Delete mark

# Amortized Analysis of DecreaseKey



decrease-key of x from 35 to 5

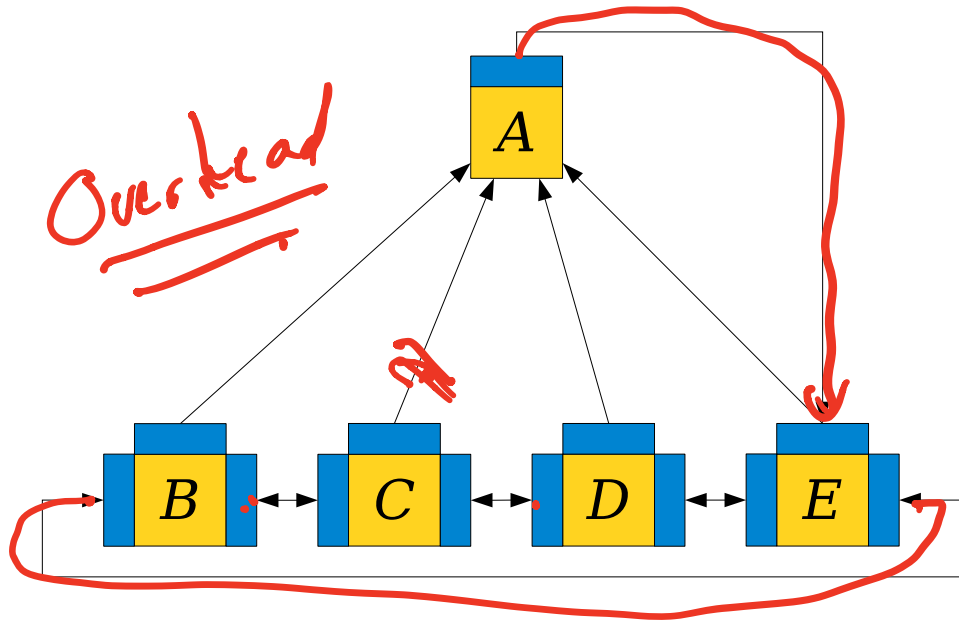second child cut

but don't mark parent if it's a root

20

# Amortized Analysis of DecreaseKey

- DecreaseKey: Actual cost O(C) where C number of cuts
  - Number of trees increases by number of cuts: potential increase
  - Number of marked nodes decreases by number of cuts - 1: May mark one new node, remove marks from others cut

- Amortized cost: $\Theta(C) + \Delta\Phi = \Theta(C) + (2 - 2C + C) = \Theta(1)$

- Note: rank($\mathscr{H}$) does not increase!  It can only decrease, and may only increase during DeleteMin.

- So, what's the problem?

# Fibonacci Heaps: Implementation

- In order to do cuts efficiently O(1), must have very complex data structures
- Children in doubly-linked circular lists
  - Point to parent
  - Parent points to one child in list
- Awful linked lists!
  - But now, can do in O(1):
    - Cut node from parent
    - Add another child to node

# Fibonacci Heaps: Implementation

- Size of Fibonacci Heap node: each node in a Fibonacci heap stores
  - A pointer to its parent.
  - A pointer to the next sibling.
  - A pointer to the previous sibling. A pointer to an arbitrary child.
  - A bit for whether it's marked.
  - Its order.
  - Its key.
  - Its element

- In practice, Fibonacci heaps are slower because of overheads
  - Good for theoretical analysis
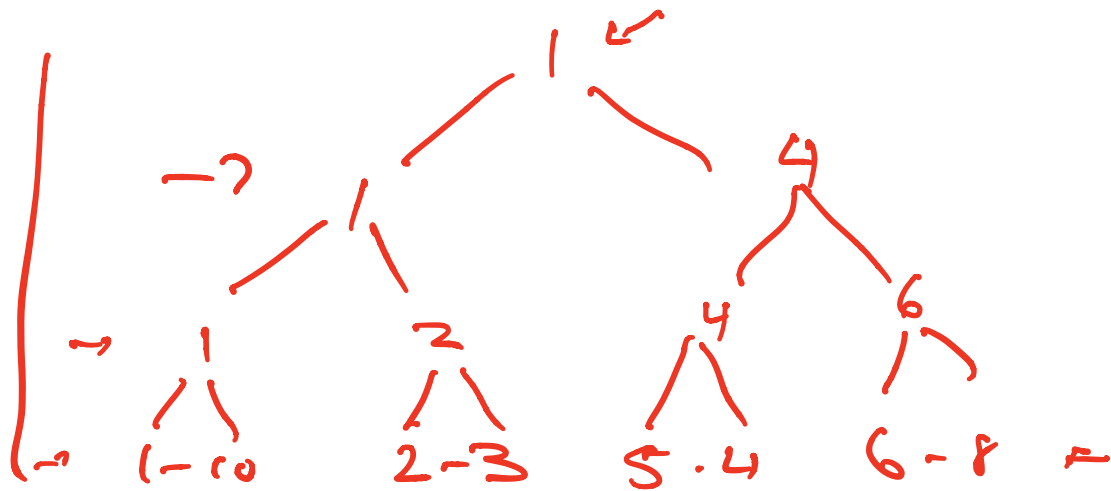  - Good research topic: Simpler heaps with O(1) DecreaseKey
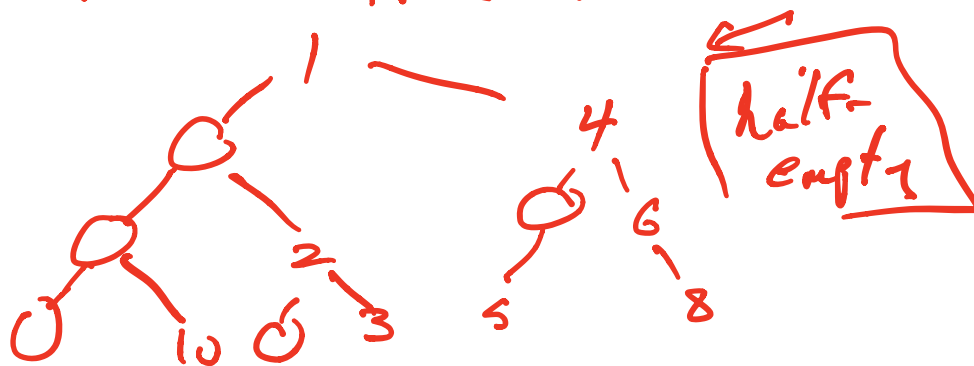
1986 - - -->

2̃3 years of search for better than...

2008-2011 → New ideas!
   Lazy merge + Simple cut...

## Quake heaps...



Tournament Tree...



half-
empty

Representation. →
   Winner

winner:



Binomial
Tree!