



A FASTER EXTERNAL MEMORY PRIORITY QUEUE WITH DECREASE KEYS

[LARSEN, JIANG] SODA19

AVIEL ZECHARIA

AUTHORS



Shunhua Jiang



Kasper Green
Larsen

AGENDA

- X Background
 - Definitions
 - The computational model
- X Applications & Motivation
- X Previous results overview
- X A faster priority queue



BACKGROUND

PRIORITY QUEUE [RAM MODEL]

- X Fibonacci heap (Strict-Fibonacci for de-amortization)
 - `DeleteMin()` – $O(\log(N))$
 - `FindMin()` – $\theta(1)$
 - `Insert(k, p)` – $\theta(1)$
 - `DecreaseKey(k, p)` – $\theta(1)$
- X Many interesting implementations
 - Tournament Tree

EXTERNAL MEMORY MODEL

- X The cost of an algorithm is the **number of memory transfers required**, operations on cached data are considered free
- X The parameters of the model are as follow:
 - N – The total size of the problem in external-memory (`HDD`)
 - M – The number of items can fit into the local memory (`RAM`)
 - B – The number of items in one block of data (`PAGE`)
- X Cache-Oblivious Model
 - M and B are not explicit known (multi-level memory hierarchy)
 - Won't be discussed in this session

FOLKLORE

- X Clearly, $T(N)$ memory accesses algorithm in the RAM model can be trivially converted into $T(N)$ memory transfers – by ignoring locality
- X Optimal Scanning
 - $O(N/B)$ by continuous scanning (B-factor improvement 😊)
- X Optimal Searching
 - $O(\log_B(N))$ using B-Trees ($\log(B)$ -factor improvement 😞)
- X Optimal Sorting
 - $O(\frac{N}{B} \log_{M/B}(\frac{N}{B}))$ using M/B merge-sort
 - computationally equivalent to priority queues (in both models)

KNOWN RESULTS

[KUMAR, SCHWABE 96]

X Priority Queue (with DecreaseKey)

- FindMin() – $\theta(1)$
- DeleteMin, Insert, DecreaseKey – $O(1/B \log_2(N/B))$ amortized
- Can we do better?



APPLICATIONS

YOUR OS !

Matrix scanning

X for i in 1 to ROW:
 for j in 1 to COL:
 ACCESS(MAT, i, j)

X for j in 1 to COL:
 for i in 1 to ROW:
 ACCESS(MAT, i, j)

1	2	3	4	B_0
5	6	7	8	B_1
9	10	11	12	B_2
13	14	15	16	B_3

YOUR FS !

Modern file-systems (Ext4/BTRFS/...)

- X Using B-tree variations for disk blocks indexing
- X Only 3-4 levels for full HDD indexing
 - X Optimal number of I/O operations

ALGORITHMS

X Dijkstra SSSP

- $O(V)$ times DeleteMin
- $O(E)$ times DecreaseKey

X Prim MST

- $O(V)$ times DeleteMin
- $O(E)$ times DecreaseKey

X etc.

3.

PREVIOUS RESULTS

AN I/O-EFFICIENT TOURNAMENT TREE

[KUMAR, SCHWABE 96]

- X “Static” tournament tree with $\theta(N/M)$ leaves $\rightarrow \theta(\log_2(N/M))$ height
- X Each leaf represent the i 'th key with initialized priority of infinity
- X Each internal node contains $O(M)$ memory:
 - ‘Winners’ buffer of size M
 - ‘Signal’ buffer of size M (order matters !)
- X The root node is always loaded in local memory
- X Invariants
 - Winners buffer (real) size must be greater than $M/2$
 - Signals buffer size must be smaller than M

DS OVERVIEW

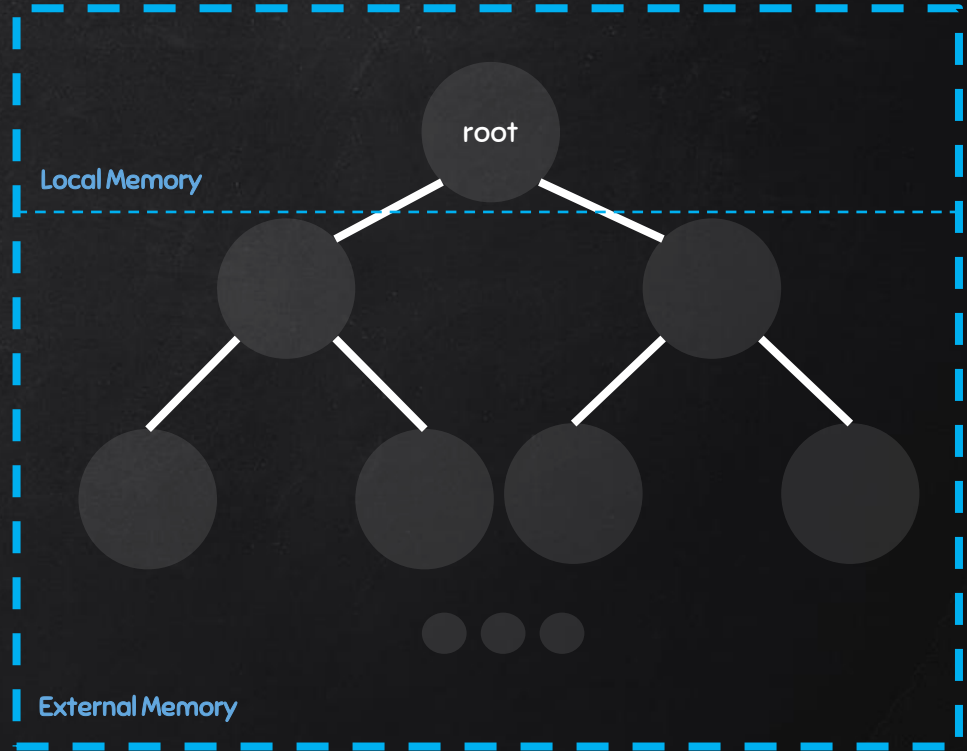


Winners

$O(M)$

(3, 5)	(2, 7)	
U(3, 5)	D(7)	U(8, 2)

Signals



UPDATE KEY WITH PRIORITY SIGNAL

Changes key's priority only if the new priority is smaller than the old one

- X If the key exists in the winners buffer
 - Try to update key priority, otherwise **drop the signal**
- X Otherwise, if there is a greater priority key'
 - Replace (maximal priority) key' with key in the winning buffer
- X Add the operation to the signals buffer
- X If the signals buffer is full
 - *Empty* procedure



ROOT

Winners

Signals

DELETE KEY (AND EXTRACTMIN) SIGNAL

Changes the priority of the key to infinity.

- X If the key exists in the winners buffer
 - Delete the key from the winners buffer
- X Add the operation to the signals buffer
- X If the signals buffer is full
 - *Empty* procedure
- X If the winners buffer is too sparse
 - *fill-up* procedure



ROOT

Winners

Signals

EMPTY PROCEDURE

Signal buffer is full and must be emptied

- X Forward signals to the children
 - Load children to memory
 - Apply & push parent signals
 - Signals order must be kept (FIFO)
- X Recursively



Winners

Signals

Winners

Signals

FILL-UP PROCEDURE

Winners buffer is too sparse and must be filled

X *Empty* procedure

- In order to force parent signals before children loading.

X Load & Update M elements from children

- Minimal priorities.
- Parent signals have already applied !

X Recursively



Winners

Signals

Winners

Signals

AMORTIZED ANALYSIS

Standard credits argument

- X Real signal logic does not cost any I/O's
 - Anyway, we will pay $O(1/B)$ for each vertex in the path to Leaf(k)
- X Empty procedure called after $O(M)$ new signals
 - Signal may only prorogate down
- X Fill-Up procedure called after $O(M)$ delete signals occurred
 - No signals are propagate up
- X On an IO-efficient tree with N elements, a sequence of k operations (Delete\DeleteMin\Update) requires at most $O(K/B \log_2(N/B))$ I/O's



4.

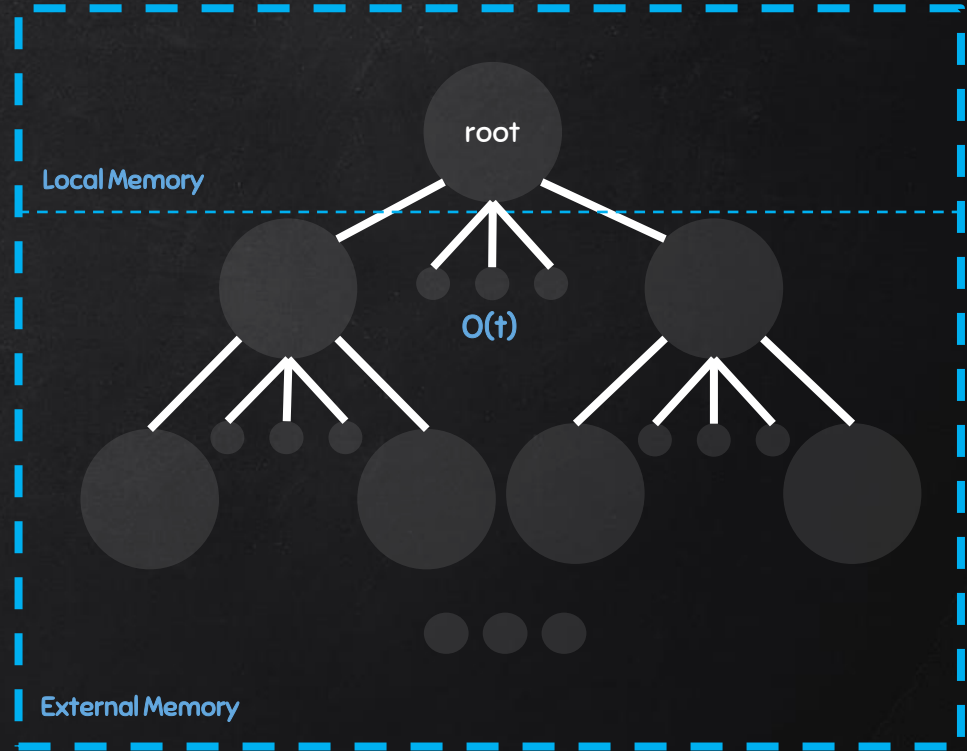
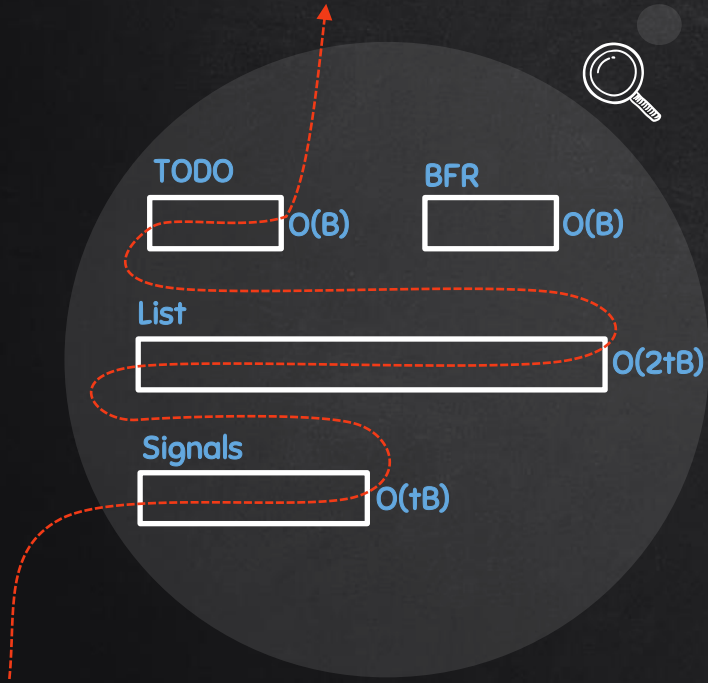
A FASTER PRIORITY QUEUE

A FASTER PRIORITY QUEUE

[LARSEN, JIANG 19]

- X “Static” t -ary tree with $\theta(N/tB)$ leaves $\rightarrow \theta(\log_2(N/tB))$ height
 - $t \sim \log^{0.01} N$
- X Each leaf represent the i 'th key with initialized priority of infinity
- X Each internal node contains $O(tB)$ memory:
 - List of elements of size $2tB$
 - Signals buffer of size tB
 - **Bloom Filter Replacement of elements of size B**
 - **TODO buffer of size B**
- X The root node is always loaded in local memory

DS OVERVIEW



BLOOM FILTER REPLACEMENT

[PAGH, RAO 05]

- X Insertion & Deletion can be done in amortized expected constant time
- X The data structure can only make false positive errors with probability at most ε
- X The space usage is at most $(1 + o(1))n \log^{1/\varepsilon} + O(n + w)$ **bits**
- X In our case
 - $n = 2tB, t = \log^{0.01} N, w = \log N, \varepsilon = 1/\log^3 N$
 - Space complexity $O(B)$ **words**

INVARIANTS

- X Each key is uniquely stored in the tree
 - Otherwise, there must exist Delete signals that will delete the redundant copies & lower update signals
- X TODO buffer may only contain single Delete | Update | Delete→Update
- X The deeper you go, the lower the priority

DELETE KEY (AND EXTRACTMIN) SIGNAL

Changes the priority of the key to infinity.

- X If there exists an entry with key k in the root
 - Delete the key from the list
- X Otherwise
 - Push Delete signal



ROOT

List

Signals

UPDATE KEY WITH PRIORITY SIGNAL

Changes key's priority only if the new priority is smaller than the old one

- X If there exists an entry (k, p') in the list of the root
 - Update to $(k, \min(p, p'))$
- X Otherwise, if $p > \text{Boundary}(\text{root})$
 - Push Update signal
- X Otherwise
 - Insert (k, p) to the list
 - Push Delete signal



ROOT

List

Signals

EMPTY-LIST PROCEDURE

Node's list has more than 2tB entries and must be restored to tB entries

X Call *apply-TODO* & *push-signal* procedures

X For each element in the list

- If child's TODO buffer contains Delete, replace it with Update
- Otherwise, append it to the children's list
- Update BFRs & Boundaries



TODO

BFR

List

Signals



TODO

BFR

List

Signals

FILL-UP PROCEDURE

Node's list is empty, and must be filled with tB minimal entries

- X Call *apply-TODO* & *push-signal* procedures
- X Iterate the children's B lowest elements (considering TODO)
 - Add the element to the node's list
 - Delete from children's list & TODO buffer.
 - Update BFRs & Boundaries



TODO

BFR

List

Signals



TODO

BFR

List

Signals

PUSH-SIGNAL PROCEDURE

Pushes the signals in the signal buffer down to its children

X Delete(k) signal

- If k is “actual” in c – insert delete to TODO, else to signal buffer

X Update(k, p) signal

- if $p < \text{Boundary}(c)$, update TODO & push delete signal
- Otherwise
 - If k is “actual” in c – update TODO & push delete signal
 - Else, push update signal



TODO

BFR

List

Signals



TODO

BFR

List

Signals

APPLY-TODO PROCEDURE

Node's TODO buffer is full and must be applied to node's list entries.

X For each TODO signal

- If false-positive update signal detected, push to signal buffer
- Else, exactly the same as Delete & Update operations !



TODO

BFR

List

Signals

I/O COMPLEXITY

Standard credits argument

- X Real signal logic does not cost any I/O's
 - Anyway, we will pay $O(1/B)$ for each vertex in the path to $\text{Leaf}(k)$
- X Push Signal – each signal goes into at most one signal buffer at each level $O(h/B)$
- X Apply TODO – In expectation, each signal goes into $1 + \epsilon h < 2$ TODO buffers $O(t/B)$

I/O COMPLEXITY

- X Empty List – using the update signal credits, each entry could empty at most one entry at each level $O(h/b)$
- X Fill Up – same has Empty List with delete signals credits
- X The expected amortized cost for each Delete, Update, ExtractMin operation is $O(h/B + t/B) = O(\frac{1}{B} \log_2 \frac{N}{B} / \log \log N)$ I/O's

FURTHER WORK

- X Cache-oblivious model
 - “Cache-Oblivious Priority Queues with Decrease-Key and Applications to Graph Algorithms” [Iacono, Jacob, Tsakalidis 20’]
- X Without DecreaseKey, there is a tight bound $\theta(1/B \log_{M/B}(N/B))$
- X There is still a gap between the algorithm complexity and the known lower bound of $O(1/B \log(B)/\log(\log(N)))$
 - The lower bound does not include the Delete(k) operation.

CONCLUSION

- X This data structure improves the external memory priority queue with DecreaseKey for the first time in over a decade
- X Our computers designed to work in the external-memory model
- X The priority queue has a direct impact on various core algorithms



THANKS!

Any questions?

Link to the paper: <https://arxiv.org/pdf/1806.07598.pdf>