madalgo - -**CENTER FOR MASSIVE DATA ALGORITHMICS**

A Tight Lower Bound for Dynamic Membership in the External Memory Model

Problems and Motivations	
 Membership: Maintain a set S in the universe U with S ≤ n. Given an x in U, answer whether x is in S or not? Dictionary: Same as Membership. Except that if x is in S, then we have to return x's associate information. Goal: Tradeoff between the (amortized) update cost t_u and the query cost t_q in the external memory (EM) model. Motivations: Two of the most fundamental data structure problems in computer science! Used extensively in databases, routers and search engines. 	Query n^{ϵ} $\log_B \log n n$ $1 \circ 0$
The Conjecture	
A long-time folklore conjecture in the external memory community (explicitly stated by Jensen and Pagh, 2007): $\mathbf{t}_{u} \text{ must be } \Omega(1) \text{ if } \mathbf{t}_{q} \text{ is required to be } O(1)$ $\mathbf{t}_{u} : \text{ expected amortized update cost.} \\ \mathbf{t}_{q} : \text{ expected average query cost.}$ Intuitively speaking, we cannot buffer the updates without sacrificing the query cost much.	 A sharp thres either use extended on the extended on t
Selected Previous Results	
 Knuth [2] analyzed <i>external hashing</i>: Expected average cost of an operation is 1 + ½ ^{Ω(B)}, provided that the load factor is less than a constant smaller than 1. (Assuming truly random hash function exists.) Arge [3] presented a data structure, <i>buffer tree</i>, supporting amortized updates cost O(1/B^{1-ε} log n) and query cost O(log_B n). Two results in two extremes. 	 [1] E. Verbin an <i>the External</i> [2] D. E. Knuth. <i>Programmin</i> [3] L. Arge. <i>The Structures</i>. <i>A</i>



Our Results



Striking Implications

shold result: for the external dictionary/membership problem, ternal hashing, or use the buffer tree, there is nothing in between!

impossible in the EM model with sublogarithmic query time.

we can show that the query complexities of many problems such reporting, predecessor, partial-sum, etc., are all the same in the the update time is less than 1.

a "clearner" model than RAM in certain perspectives.

References

nd Q. Zhang. A Tight Lower Bound for Dynamic Membership in Memory Model. STOC 2010.

Sorting and Searching. Volume 3 of The Art of Computer ng. Addison-Wesley, Reading MA, second edition, 1998.

Buffer Tree: A Technique for Designing Batched External Data Algorithmica 2003.

Framework:

Insert n random items. Pick a random snapshot "END" in the insertion sequence, and divide previous insertions into blocks with exponentially increasing size. We find a set of cells C_i for each block. C_i are disjoint. Finally, we argue that a random query needs to read one bit from at least $\Omega(\log_{B \log n} n)$ blocks.



Proof ideas:

- arguments.
- paths intersecting C_i. We are done.



MADALGO – Center for Massive Data Algorithmics, a Center of the Danish National Research Foundation



Technical Contributions

• For each block i, we first prove that for 0.99 fraction of items inserted in block i, their query paths will intersect C_i. The proof is by encoding

Next, we introduce a lemma called LOSI, and then use it to prove that if 0.99 fraction of newly inserted items having query paths intersecting C_i , then at least a constant fraction of items in the universe have query