



ANIL SHANBHAG

MIT

J/W Alekh Jindal, Sam Madden, Jorge Quiane,Aaron J. ElmoreMicrosoftMITQCRIUniv. Chicago

Today

Data collection is cheap => Lots of data !



Data Partitioning

Find average order size for all orders between Sept 10 and Sept 11, 2017



Data Skipping - Skip data blocks not necessary

10% selectivity query => 10x faster if data partitioned on selection predicate

The Problem



How to get benefits of partitioning in this case ?

Our Approach

Do everything adaptively !

Two step process:

- 1. Upfront load the dataset partitioned
- 2. As users query, *incrementally* improve the partitioning of the data



Distributed storage systems like HDFS, files broken into blocks (128 MB chunks)

Upfront Partitioning

> Instead of partitioning by size, partition by attributes.

> Same number of blocks created as in HDFS. Each block now has additional metadata





Adaptive Re-Partitioning

When user submits a query, optimizer tries to improve the partitioning by reorganizing the partitioning tree

Here if queries ask A <= 3 many times, replace B₇ by A₃

Done on datasets which are O(1TB) with ~ 8000 node partition trees.

System Architecture



Predicated Scan Query Example:

FIND employees WITH Age < 30 AND 20k < Salary < 40k

1. Upfront Partitioner

Goal: Generate a partitioning tree WITHOUT an upfront query workload

> Generates a tree with heterogeneous branching

> Balance the partitioning benefit across all attributes



Allocation

Goal: Balance partitioning benefit across attributes

Allocation of attribute i ~ average partitioning of an attribute j

$$=\Sigma_{\text{all nodes i}}$$
 nij Cij







2. Adaptive Query Executor

Goal: Return matching tuples + <u>check if partitioning layout can be improved</u>

Alternatives found via transformations on the partitioning tree



Getting a plan

Input : Node node, Predicate pred
1 if isLeaf(node) then
2 return $P_{none}(node);$
3 else
4 if isLeftAccessed(node) then
5 [leftPlan \leftarrow getSubtreePlan(node.lChild, pred);
6 if isRightAccessed(node) then
7 rightPlan \leftarrow getSubtreePlan(node.rChild, pred)
/* consider swap */
8 if leftPlan.fullyAccessed and rightPlan.fullyAccessed then
9 currentCost $\leftarrow \sum_i \text{Cost}(node, q_i);$
10 what If Node \leftarrow clone(<i>node</i>);
11 whatIfNode.predicate \leftarrow newPred;
12 swapNode(<i>node</i> , <i>whatIfNode</i>);
13 newCost $\leftarrow \sum_i \text{Cost}(whatIfNode, q_i);$
14 benefit = currentCost - newCost;
15 if benefit > 0 then
16 updatePlanIfBetter(<i>node</i> , $P_{swap}(node, whatIfNode)$);
/* consider pushup */
17 if leftPlan.ptop and rightPlan.ptop then
<pre>18 updatePlanIfBetter(node, P_{pushup}(node, node.lChild, node.rChild));</pre>
/* consider rotate */
19 if node.attribute == predicate.attribute then
20 if leftPlan.ptop then
21 updatePlanIfBetter(node, P _{rotate} (node, node.lChild));
22 if rightPlan.ptop then
23 updatePlanIfBetter(node, Protate(node, node.rChild));
/* consider doing nothing */
<pre>updatePlanIfBetter(node, P_{none}(node));</pre>
25 return node.plan;

Algorithm 3: getBestPlan **Input** : Tree tree, Predicate[] predicates 1 while predicates $\neq \phi$ do prevPlan \leftarrow tree.plan; 2 **foreach** *p* in predicates **do** 3 Plan newPlan \leftarrow getSubtreePlan(*tree.root*, *p*); 4 updatePlan(tree.root, newPlan, p); 5 **if** *tree.plan* \neq *prevPlan* **then** 6 remove from predicates the newly inserted predicate; 7 else 8 break; 9 10 **if** *tree.plan.benefit* > *tree.plan.cost* **then return** tree.plan; 11 12 **else** return null; 13

Cost Model

The system maintain a window W of past queries

Compute Benefit and Repartitioning Cost for the best plan

Repartitioning **ONLY** happens when reduction in the total cost of the query workload is greater than re-partitioning cost.

Solves constant re-partitioning due to random query sequences and bounds the worse case impact.

$$\operatorname{Cost}(T,q) = \sum_{b \in \operatorname{lookup}(T,q)} n_b$$

Benefit
$$(T') = \sum_{q \in W} Cost(T, q) - \sum_{q \in W} Cost(T', q)$$

RepartitioningCost
$$(T,q) = \sum_{b \in B} c \cdot n_b$$

Performance

4 metrics

- 1) Load time
- 2) Time taken by first query
- 3) Aggregate runtime over a workload
- 4) Incremental improvement with workload hints

Load Time

TPC-H: Scale Factor 200 + De-normalized. Data size:1.4TB



Loading performance: 1.38 times slower than HDFS

Load time scales almost linearly with data size and independent of number of columns

Time taken by first query



On Average: 45% better than full scan 20% better than k-d tree

Aggregate Workload Runtime



Workload: 200 Queries generated from random initialization of 8 query templates of TPC-H benchmark

full scan – Baseline

range – partitions on orderdate (1 per date) 1.88x better

range2 - partitions on orderdate(64), r_name(4),c_mktsegment(4),quantity(8) 3.48x better



Workload Hints



Better Init: Starts with custom allocation to mimic *range2*

6.67x better than fullscan

Filtering ratio: default : 0.81 better init : 0.9

Conclusion

•Amoeba is a distributed storage system based on an adaptive data partitioning scheme

- Low loading overhead
- Improved first query performance
- Adapt to changes and significantly improvement to workload runtime
- Can exploit workload hints

•Allows analysts to get started right away and reap benefits of partitioning without an upfront workload