

Relax and Let the Database do the Partitioning Online

Alekh Jindal and Jens Dittrich

Information Systems Group, Saarland University
<http://infosys.cs.uni-saarland.de>

Abstract. Vertical and Horizontal partitions allow database administrators (DBAs) to considerably improve the performance of business intelligence applications. However, finding and defining suitable horizontal and vertical partitions is a daunting task even for experienced DBAs. This is because the DBA has to understand the physical query execution plans for each query in the workload very well to make appropriate design decisions. To facilitate this process several algorithms and advisory tools have been developed over the past years. These tools, however, still keep the DBA in the loop. This means, the physical design cannot be changed without human intervention. This is problematic in situations where a skilled DBA is either not available or the workload changes over time, e.g. due to new DB applications, changed hardware, an increasing dataset size, or bursts in the query workload. In this paper, we present AUTOSTORE: a self-tuning data store which rather than keeping the DBA in the loop, monitors the current workload and partitions the data automatically at checkpoint time intervals — without human intervention. This allows AUTOSTORE to gradually adapt the partitions to best fit the observed query workload. In contrast to previous work, we express partitioning as a One-Dimensional Partitioning Problem (1DPP), with Horizontal (HPP) and Vertical Partitioning Problem (VPP) being just two variants of it. We provide an efficient O^2P (One-dimensional Online Partitioning) algorithm to solve 1DPP. O^2P is faster than the specialized affinity-based VPP algorithm by more than two orders of magnitude, and yet it does not lose much on partitioning quality. AUTOSTORE is a part of the OctopusDB vision of a One Size Fits All Database System [13]. Our experimental results on TPC-H datasets show that AUTOSTORE outperforms row and column layouts by up to a factor of 2.

Key words: changing workload, online partitioning, self-tuning

1 Introduction

Physical database designs have been researched heavily in the past [23, 3, 2, 5, 6, 31, 7, 22, 26, 12]. As a consequence, nowadays, most DBMSs offer design advisory tools [1, 32, 30, 4]. These tools help DBAs in defining indexes, e.g. [8], as well as horizontal and/or vertical partitions, e.g. [3, 15]. The idea of these tools is to analyze the workload at a given point in time and suggest different

physical designs. These suggestions are computed by a what-if analysis. What-if analysis explores the possible physical designs. However, just finding the right set of partitions is NP-hard [29]. Therefore the search space must be pruned using suitable heuristics, i.e. typically some greedy-strategy [2]. The cost of each candidate configuration is then estimated using an existing cost-based optimizer, i.e. the optimizer is tricked into believing that the candidate configuration already exists. Eventually, a suitable partitioning strategy is proposed to the DBA who then has to re-partition the existing database accordingly.

1.1 Problems with Offline Partitioning

The biggest problem with this approach is that it is an *offline* process. The DBA will only reconsider the current physical design at certain points in time. This is problematic. Assume the workload changes over time, e.g. changes in the workload due to new database applications, an increasing dataset size, or an increasing number of queries. In these situations the existing partitioning strategies should be revisited to improve query times. In the current *offline approach* however, the partitioning strategies will only be changed if a human — the DBA — triggers an advisory tool with the most recent query logs and eventually decides to repartition the data. This means, the vertical and horizontal partitioning strategies are carved in stone until the DBA changes them eventually. Furthermore, current advisory tools attempt to find near optimal partitioning strategies, which is very expensive. This is especially problematic if the database system has to handle bursts and peaks. For instance consider (i) a ticket system selling 10 million Rolling Stones tickets within three days; (ii) an online store such as Amazon selling much higher volumes before christmas; or (iii) an OLAP system having to cope with new query patterns. In these types of applications it is not acceptable for users to wait for the DBA and the advisory tool to reconfigure the system. If the system stalls due to a peak workload, the application provider may lose a lot of money.

1.2 Research Goals and Challenges

Our **goal** is to research a database store that decides on the suitable partitioning strategy *automatically*, i.e. without any human intervention. As the search space of possible partitions is huge [29], it is clear from the beginning that an *optimal* automatic partitioning is not always feasible. However, we believe that an automatic (or: *online*) partitioning will in most cases be much better than the one suggested by even a skilled DBA — similarly the physical query execution plans being in most cases better than hand-crafted plans. The risk of not reaching optimality is similar to the risk of adaptive indexing [14] and cracking [19]. However, the possible gains of such an approach may be similarly tremendous.

This leads to interesting **research challenges**:

- (1.) There exist a plethora of state-of-the-art offline algorithms, e.g. [23, 24, 2], for suggesting suitable vertical and horizontal partitions. However, given the

huge search space, it turns out that their runtime complexity is unacceptably high to be applicable in an online setting where the available time to decide on a new partitioning is rather limited. Therefore we must develop new algorithms.

- (2.) We need algorithms that decide on data partitioning automatically and while the database is running, i.e. take decisions to repartition the data.
- (3.) Any automatic repartitioning must not block or stall the database and/or access to entire tables, a problem more likely in archival disk-based databases.

1.3 Contributions

In this paper, we present AUTOSTORE, a fully automatic database store, to solve these challenges. To the best of our knowledge, this work is the first to solve the database partitioning problem with a fully automatic *online approach*. The contributions of this paper are:

- (1.) We express partitioning as general One-Dimensional Partitioning Problem (1DPP), with Vertical (VPP) and the Horizontal Partitioning Problem (HPP) as subproblems of it. Both subproblems may be solved by solving 1DPP (Section 2).
- (2.) We present AUTOSTORE, an online self-tuned database store that is a step towards implementing the OctopusDB vision [13, 21]. The core components of AUTOSTORE are: dynamic workload monitor, partitioning unit clusterer, partitioning analyzer, and partitioning optimizer (Section 3).
- (3.) We present an online database partitioning algorithm O²P (One-dimensional Online Partitioning) to solve 1DPP (Section 4).
- (4.) We show an extensive evaluation of our algorithm over TPC-H and SSB benchmarks. We present experimental results from mixed OLTP/OLAP workloads over a main-memory and a BerkeleyDB implementation (Section 5).

2 Vertical and Horizontal Partitioning

2.1 Preliminaries

Typically, users partition databases horizontally based on data value ranges, hashes, or lists. This is because data values are comparable across a column. However, this is not true for data values across a row. Therefore, sophisticated partitioning methods have been developed for VPP. In this section we will revisit the basics of VPP. This also serves as the ground work for our 1DPP.

Naïve Approach. Database researchers pointed out the heuristic [16] and NP-hard [29] nature of partitioning problem pretty early. The number of ways to partition vertically, for x attributes, is given by bell number $B(x)$. The naïve approach to find the optimal solution is to enumerate all bell numbers. The complexity of the naïve approach is $O(x^x)$, making it infeasible for large databases.

Affinity based Approach. The naïve approach considers all possible partitions, even the ones having attributes which are *never* accessed together. To address this, *attribute affinity* was introduced as a measure of pairwise attribute

similarity [18, 23, 10, 24, 11]. The core idea of affinity based partitioning is to compute affinities between every pair of attributes and then to cluster them such that high affinity pairs are as close in neighborhood as possible. To compute affinity between different attributes, we need to know their access patterns. A *usage function* $U(q, a)$ denotes whether or not query q references attribute a . $U(q, a) = 1$ if q references a and 0 otherwise. For example, in TPC-H Lineitem table, $U(Q_1, \text{PartKey}) = 1$ as Query 1 references attribute PartKey. The usage function may also be extended to incorporate query weights, reflecting the importance levels or relative frequencies of queries. To measure the affinity between two attributes a_i and a_j , the *affinity function* $A(a_i, a_j)$ simply counts their co-occurrences in the query workload, i.e. $A(a_i, a_j) = \sum_q U(q, a_i) \cdot U(q, a_j)$. For instance, in Lineitem table, $A(\text{PartKey}, \text{SuppKey}) = 5$, as attributes PartKey and SuppKey co-occur in five queries. The affinity function produces a 2D affinity matrix between every pair of attributes. The goal now is to cluster the matrix such that the cells having similar affinity values are placed close together in the matrix. Every order of rows and columns in the matrix gives a new *ordering of attributes* (\preceq). For example, consider the following affinity matrices for PartKey, SuppKey, and Quantity attributes in TPC-H Lineitem table.

The left matrix represents an attribute ordering $\text{PartKey} \preceq \text{SuppKey} \preceq \text{Quantity}$,

	PartKey	SuppKey	Quantity
PartKey	8	5	6
SuppKey	5	8	4
Quantity	6	4	9

whereas the right matrix represents ordering $\text{PartKey} \preceq \text{Quantity} \preceq \text{SuppKey}$.

	PartKey	Quantity	SuppKey
PartKey	8	6	5
Quantity	6	9	4
SuppKey	5	4	8

Given attribute ordering \preceq , an *affinity measure* $M(\preceq)$ measures the quality of the affinity clustering as $M(\preceq) = \sum_{i=1}^x \sum_{j=1}^x A(a_i, a_j)[A(a_i, a_{j-1}) + A(a_i, a_{j+1})]$. It holds that $A(a_0, a_j) = A(a_i, a_0) = A(a_{x+1}, a_j) = A(a_i, a_{x+1}) = 0$. For the left matrix above $M(\preceq) = 404$ and for the right matrix $M(\preceq) = 440$. Indeed, the right matrix has better clustering since affinity between attributes PartKey and Quantity (=6) is more than that between PartKey and SuppKey (=5). Thus, the objective of affinity matrix clustering problem now is to maximize the affinity measure. One (greedy) approach is to place the attributes one-by-one such that the *contribution to the affinity measure* at each step is maximized [23]. The contribution to the affinity measure of a new attribute a_k when placed between two already placed attributes a_i and a_j is: $Cont(a_i, a_j, a_k) = 2 \cdot \sum_{z=1}^n [A(a_z, a_i) \cdot A(a_z, a_k) + A(a_z, a_k) \cdot A(a_z, a_j) - A(a_z, a_i) \cdot A(a_z, a_j)]$. In this clustering approach, we first place two random attributes; then, in the neighborhood, we place the attribute which maximizes the contribution to the affinity measure. We repeat this process until all attributes are placed.

2.2 Problem Statement

In this section we express HPP and VPP as a general 1DPP. The first step to do so is to identify the smallest indivisible units of storage.

Definition 1. A *partitioning unit set* $P_u = \{u_1, u_2, \dots, u_n\}$ is the set of n smallest pieces of data.

Definition 2. A *partitioning unit ordering* \preceq defines an order on the partitioning units in P_u .

Partitioning units could be attributes along the vertical axis or tuples along the horizontal axis. However, partitioning at the tuple level may not make sense due to large number of partitioning units and hence high complexity. Therefore, we usually consider sets of tuples, based on some key, as partitioning units (horizontal partitioning). Similarly, we could also consider groups of columns as partitioning units (vertical partitioning). Below, we introduce some new concepts needed for our one-dimensional partitioning problem statement. First, we express partitioning as a logical partitioning, to be able to use it in an algorithm.

Definition 3. A *split vector* S is a row vector of $(n-1)$ split lines in ordering \preceq , where a split line s_j is defined between partitioning units u_j and u_{j+1} as follows:

$$s_j = \begin{cases} 1 & \text{if there is split between } u_j \text{ and } u_{j+1} \\ 0 & \text{for no split} \end{cases} .$$

A split vector S captures the logical partitioning over a given dataset. For instance, a split vector $S_1=[0,0,0,1,0,1,1]$ corresponds to a partitioning of $u_1, u_2, u_3, u_4|u_5, u_6|u_7|u_8$. However, in order to estimate costs using a cost-based query optimizer, a split vector still needs to be translated in terms of partitioning units:

Definition 4. A *partition* $p_{m,r}(S, \preceq)$ is a maximal chunk of adjacent partitioning units from u_m to u_r , such that split lines s_m to s_{r-1} are all 0.

Definition 5. A *partitioning scheme* $P(S, \preceq)$ over relation R is a set of disjoint and complete partitions, i.e.

$$\begin{aligned} \bigcup_x p_{m_x, r_x}(S, \preceq) &= R, \\ p_{m_x, r_x}(S, \preceq) \cap p_{m_y, r_y}(S, \preceq) &= \phi, \forall x, y \text{ such that } x \neq y. \end{aligned}$$

Partitioning scheme expresses the actual arrangement of partitioning units, given a split vector. For instance, for split vector S_1 , partition $p_{1,4}(S_1, \preceq)$ is $\{u_1, u_2, u_3, u_4\}$ and partitioning scheme $P(S_1, \preceq) = \{p_{1,4}(S_1, \preceq), p_{5,6}(S_1, \preceq), p_7(S_1, \preceq), p_8(S_1, \preceq)\}$. Finally, in order to evaluate partitioning schemes in an online setting, we need to model the online query workload.

Definition 6. An *Online Workload* W_{t_k} is a stream of queries $\{q_0, \dots, q_{t_k-1}, q_{t_k}\}$ seen till time t_k , where $t_k > t_{k-1} > \dots > 0$.

Further, let $C_{\text{est.}}(W_{t_k}, P(S, \preceq))$ denote the execution cost of workload W_{t_k} as estimated by a cost-based optimizer. Now, we express our one-dimensional partitioning problem as follows.

One-dimensional Online Partitioning Problem. Given an online workload W_{t_k} and partitioning unit ordering \preceq , find the split vector S' that minimizes the estimated workload execution cost, i.e.

$$S' = \underset{S}{\operatorname{argmin}} C_{\text{est.}}(W_{t_k}, P(S, \preceq)). \quad (1)$$

The complexity of the above problem depends on the number of partitioning schemes P , which in turn depends on the range of values of split vector S . Note that the one-dimensional partitioning problem has the following sub-problems: (1) *Vertical Partitioning*, if the partitioning unit set P_u is a set of attributes, and (2) *Horizontal Partitioning*, if the partitioning unit set is a set of horizontal ranges.

In the next section we describe our AUTOSTORE system and discuss how it solves the one-dimensional partitioning problem in an online setting.

3 AutoStore

In this section we present AUTOSTORE, an automatically and online partitioned database store. The workflow in AUTOSTORE is as follows: (1) dynamically monitor the query workload and dynamically cluster the partitioning units; (2) analyze the affinity matrix at regular intervals and decide whether or not to create the partitions; and (3) keep data access unaffected while monitoring workload and analyzing partitions. The salient features of AUTOSTORE are as follows:

- (1.) *Physical Data Independence.* Instead of exposing physical data partitioning to the user, AUTOSTORE hides these details. This avoids mixing the logical and physical schema. Thus, AUTOSTORE offers better physical data independence, which is not the case in traditional databases.
- (2.) *DBA-oblivious Tuning.* The DBA is not involved in triggering the right data partitioning. AUTOSTORE self-tunes its data.
- (3.) *Workload Adaptability.* AUTOSTORE monitors the *changes* in query workload and automatically adapts data partitioning to it.
- (4.) *Generalized Partitioning.* AUTOSTORE treats partitioning as a 1DPP. Sub-problems VPP and HPP are handled equivalently by rotating the table through ninety degrees, i.e. changing the partitioning units from attributes to ranges.
- (5.) *Cost, Benefit Optimization.* To decide whether or not to actually partition data, AUTOSTORE considers both the costs as well as the benefits of partitioning.
- (6.) *Uninterrupted Query Processing.* AUTOSTORE makes use of our online algorithm O²P which amortizes the computationally expensive partitioning analysis over several queries. The query processing remains uninterrupted.

Below we discuss four crucial components — workload monitor, partitioning unit clusterer, partitioning analyzer, and partitioning optimizer — which make online self-tuning possible in AUTOSTORE.

Workload Monitor. Online partitioning faces the challenge of creating good partitions for future queries based on the seen ones. One might consider partitioning after every incoming query. However, not only could this be expensive (due to data shuffling), the next incoming query could be entirely different from the previous one. Hence, we maintain a *query window* to capture the workload pattern and have greater confidence over partitioning decisions. Additionally, we slide the query window once it grows to a maximum size, to capture the latest workload trends. We denote a sliding query workload having N queries as

$W_{t_k}^N \subseteq W_{t_k}$. After every `CheckpointSize` number of new queries, `AUTOSTORE` triggers partitioning analysis, i.e., it takes a snapshot of the current query window and the partitioning unit ordering and determines the new partitioning.

Partitioning Unit Clusterer. The partitioning unit clusterer is responsible for re-clustering the affinity matrix after each incoming query. The affinity matrix clustering algorithm (in Section 2.1) has the following issues: (i) it recomputes all affinity values, and (ii) it reclusters the partitioning units from scratch. We need to adapt it for online partitioning in `AUTOSTORE`. The core idea is to compute all affinities once and then for each incoming query update only the affinities between referenced partitioning units. Note that the change in each of these affinity values will be 1, due to co-occurrence in the incoming query. For example, consider TPC-H Lineitem table having the affinity matrix shown at left below.

Now, for an incoming query accessing PartKey and SuppKey, only the affinities between

	PartKey	Quantity	SuppKey
PartKey	8	6	5
Quantity	6	9	4
SuppKey	5	4	8

	PartKey	Quantity	SuppKey
PartKey	9	6	6
Quantity	6	9	4
SuppKey	6	4	9

them are updated (gray cells in the right affinity matrix above). Likewise, we need to re-cluster only the referenced partitioning units. To do this, we keep the first referenced partitioning unit at its original position, and for the i^{th} referenced unit we consider the left and right positions of the $(i - 1)$ referenced units already placed. We calculate the net contribution of i^{th} referenced unit to the global affinity measure as: (*Cont* at the new position) - (*Cont* at the current position). We choose the position that offers maximum net contribution to the global affinity measure and repeat the process for all referenced partitioning units. To illustrate, in the right affinity matrix above, we first place PartKey and then consider placing SuppKey to the left (net contribution=48) and right (net contribution=0) of Partkey. Thus, we will place SuppKey to the left of PartKey.

Partitioning Analyzer. The partitioning analyzer of `AUTOSTORE` analyzes partitioning every time `CheckpointSize` number of queries are added by the workload monitor. The job of the partitioning analyzer is to take a snapshot of the query window as input, enumerate and analyze the partitioning candidates, and emit the best partitioning as output. In the brute force enumeration approach, we consider all possible values (0 or 1), for each split line in the split vector S of Equation 1. We then pick the split vector which produces the lowest estimated workload execution cost $C_{est.}(W_{t_k}^N, P(S, \preceq))$. Each split vector gives rise to a different candidate partitioning scheme. The size of the set of candidate partitioning schemes is 2^{n-1} . In Section 4 we show how the O^2P algorithm significantly improves partitioning analysis in an online setting.

Partitioning Optimizer. Given the partitioning scheme P' produced by the partitioning analyzer, the partitioning optimizer decides whether or not to transform the current partitioning scheme P to P' . The partitioning optimizer considers the expected costs of transforming the partitioning scheme as well as the expected benefits from it. We discuss these considerations below.

Cost Model. We use a cost model for full table and index scan operations over one-dimensional partitioned tables. To calculate the partitioning costs, we first

find the partitions in P which are no longer present in P' : $P_{\text{diff}} = P \setminus P'$. Now, to transform from P to P' we simply have to read each of the partitions in P_{diff} and store it back in the required partitions in P' . For instance, the transformation cost for vertical partitioning can be estimated as twice the scan costs of partitions in P . From such a transformation cost model, the worst case transformation cost is equal to twice the full table scan, whereas the best case transformation cost is twice the scan cost of the smallest partitioning unit.

Benefit Model. Same as we compute the cost of partitioning, we also need the benefit of partitioning in order to make a decision. We model partitioning benefit as the difference in the cost of executing the query window on the current and the new partitioning, i.e. $B_{\text{transform}} = C_{\text{est.}}(W_{t_k}^N, P(S, \preceq)) - C_{\text{est.}}(W_{t_k}^N, P'(S, \preceq))$.

Partitioning Decision. For each transformation made, we would have recurring benefits over all similar workloads. Hence, we need to factor in the expected frequency of the query window. This could be either explicitly provided by the user or modeled by the system. For instance, an exponential decaying model with shape parameter y : $\text{Workload Frequency}(f) = \frac{1}{1-y^{\text{MaxWindowSize}}}$ gives higher frequency to smaller query windows. AUTOSTORE creates the new partitioning only if the total recurring partitioning benefit ($\text{pBenefit} = f \cdot B_{\text{transform}}$) is expected to be higher than the partitioning cost ($\text{pCost} = C_{\text{transform}}$).

Partitioning Transformation/Repartitioning. Repartitioning data from P to P' poses interesting algorithmic research challenges. As stated before the overall goal should be to minimize transformation costs. In addition, the database or even single tables must not be stalled, i.e. by halting incoming queries. Fortunately, these problems may be solved. In a *read-only* system any table or horizontal partition may be transformed in the background, i.e. we transform P to P' and route all incoming queries to P . Only if the transformation is finished, we atomically switch to P' . For *updates*, this process can be enriched by keeping a differential file or log L of the updates that are arriving while the transformation is running. Any incoming query may then be computed by considering P and L . If the transformation is finished, we may eventually decide to merge P' with P . The right strategy for doing this is not necessarily to merge immediately. A similar discussion as for LSM trees [25] and exponential files [20] applies. For repartitioning vertical layouts there are other interesting challenges. None of them would require us to stall incoming queries or halt the database. We are planning to evaluate these algorithms in a separate study.

4 O²P Algorithm

The partitioning analyzer in Section 3 described the brute force approach of enumerating all possible values of split vector S in the one-dimensional partitioning problem. This approach has exponential complexity and hence is not desirable in an online setting. In this section, we present an online algorithm O²P (One-dimensional Online Partitioning) which solves 1DPP in an online setting. O²P does not produce the optimal partitioning solution. Instead, it uses a number of techniques to come up with greedy solution. The greedy solution is

not only dynamically adapted to workload changes, it also does not lose much on partitioning quality as well. Below we highlight the major features in O²P:

(1.) Partitioning Unit Pruning. Several partitioning units are never referenced by any of the queries in the query window. For instance, `RetailPrice` is not referenced in TPC-H `Part` table. Due to one-dimensional clustering, such partitioning units are expected to be in the beginning or the end of the partitioning unit ordering. Therefore, O²P prunes them into a separate partition right away. As an example, consider 3 leading and 2 trailing partitioning units, out of total 10 partitioning units, to be not referenced. Then, the following split lines are determined: $s_1 = s_2 = 0$, $s_3 = 1$, $s_9 = 0$, $s_8 = 1$.

(2.) Greedy Split Lines. Instead of enumerating over all possible split line combinations — as in the brute force — O²P greedily sets the best possible split line, one at a time. O²P starts with a split vector having all split lines as 0 and at each iteration it sets (to 1) only one split line. To determine which split line to set, O²P considers all split lines unset so far, and picks the one giving the lowest workload execution cost, i.e. the $(i + 1)^{th}$ split line to be set is given by: $s_{i+1} = \operatorname{argmin}_{s \in \text{unset}(S_i)} C_{\text{est.}}(W_{t_k}^N, P(S_i + U(s), \preceq))$, where $U(s)$ is a unit vector having only split line s as set; corresponding split vector is: $S_{i+1} = S_i + U(s)$.

(3.) Dynamic Programming. Observe that the partitions not affected in the previous iteration of greedy splitting will have the same best split line in the current iteration. For example, consider an ordering of partitioning units with binary partitioning: $u_1, u_2, u_3, u_4 | u_5, u_6, u_7, u_8$. The corresponding split vector is: $[0, 0, 0, 1, 0, 0, 0]$ with only split line s_4 set to 1 and all other split lines set to 0. Now, we consider all unset split lines for partitioning. Suppose s_2 and s_6 are the best split lines in the left and right partitions respectively and amongst them s_2 is the better split line. In next iteration, we already know that s_6 is the best split line in the right partition and only need to evaluate s_1 and s_3 in the left partition again. To exploit this O²P maintains the best split line in each partition and reevaluates split lines only in partitions which are further split. Since it performs only one split at a time (greedy), it only needs to reevaluate the split lines in the most recently split partition. Algorithm 1 shows the dynamic programming based enumeration in O²P. First, O²P finds the best split line and its corresponding cost in: the left and right parts of the last partition, and all previous partitions (Lines 1–6). If no valid split line is found then O²P returns (Lines 7–9). Otherwise, it compares these three split lines (Line 10), chooses the one having lowest costs (Lines 11–35), and repeats the process (Line 36).

Theorem 1. O²P produces the correct greedy result.

Theorem 2. If consecutive splits reduce the partitioning units by z elements, then the number of iterations in O²P is $\left\lceil \frac{(n-3)(n-z-1)}{2z} + 2n - 3 \right\rceil$.

Lemma 1. Worst case complexity of O²P is $O(n^2)$.

Lemma 2. Best case complexity of O²P is $O(n)$ ¹.

¹ All proofs at: <http://infosys.cs.uni-saarland.de/publications/TRAutoStore.pdf>

Algorithm 1: dynamicEnumerate

```

Input : S, left, right, PrevPartitions
Output: Enumerate over possible split vectors

1 SplitLine sLeft = BestSplitLine(S, left);
2 Cost minCostLeft = BestSplitLineCost(S, left);
3 SplitLine sRight = BestSplitLine(S, right);
4 Cost minCostRight = BestSplitLineCost(S, right);
5 SplitLine sPrev = BestSplitLine(S, PrevPartitions);
6 Cost minCostPrev = BestSplitLineCost(S, PrevPartitions);
7 if invalid(sLeft) and invalid(sRight) and invalid(sPrev) then
8   | return;
9 end
10 Cost minCost = min(minCostLeft, minCostRight, minCostPrev);
11 if minCost == minCostLeft then
12   | SetSplitLine(S, sLeft);
13   | if sRight > 0 then
14     | | AddPartition(right, sRight, minCostRight);
15     | end
16     | right = sLeft+1;
17 else if minCost == minCostRight then
18   | SetSplitLine(S, sRight);
19   | if sLeft > 0 then
20     | | AddPartition(left, sLeft, minCostLeft);
21     | end
22     | left = right;
23     | right = sRight+1;
24 else
25   | SetSplitLine(S, sPrev);
26   | if sRight > 0 then
27     | | AddPartition(right, sRight, minCostRight);
28     | end
29   | if sLeft > 0 then
30     | | AddPartition(left, sLeft, minCostLeft);
31     | end
32   | RemovePartition(sPrev);
33   | left = pPrev.start();
34   | right = sPrev+1;
35 end
36 dynamicEnumerate(S, left, right, PrevPartitions);

```

(4.) **Amortized Partitioning Analysis.** O^2P computes the partitioning lazily over the course of several queries, i.e. it performs a subset of iterations each time `AUTOSTORE` triggers the partitioning analyzer. Thus, O^2P amortizes the cost of computing the partitioning scheme over several queries. This makes sense, because otherwise we may end up spending a large number of CPU cycles, and blocking query execution, even though the partitioning may not be actually done (due to cost-benefit considerations). The partitioning analyzer returns the best split vector only when all iterations in the current analysis are done.

(5.) **Multi-threaded Analysis.** Since our partitioning analysis works on a window snapshot of the workload, O^2P can also delegate it to a separate secondary thread while the normal query processing continues in the primary thread. This approach completely separates query processing from partitioning analysis.

5 Experiments

The goal of our experiments is three-fold: (1) to evaluate the partitioning analysis in AUTOSTORE on realistic TPC-H and SSB workloads, (2) to compare the query performance of a main-memory based implementation of AUTOSTORE with No and Full Vertical Partitioning, and (3) to evaluate the performance of AUTOSTORE on a real system: BerkeleyDB. We present each of these in the following. All experiments were executed on a large computing node having Intel Xeon 2.4GHz CPU with 64GB of main memory, and running on Ubuntu 10.10 operating system.

The algorithms of Navathe et. al. [23] and Hankins et. al. [17] have similar complexity. Therefore, we label them as NV/HC. To compare and obtain a cost analysis of different components in O²P, we switch them on incrementally. Thus, we have five different variants of O²P: (i) only partitioning unit pruning (O²Pp), (ii) pruning+greedy (O²Ppg), (iii) pruning+greedy+dynamic (O²Ppgd), (iv) pruning+greedy+dynamic+amortized (O²Ppgda), and (v) pruning+greedy+dynamic+multi-threaded (O²Ppgdm).

5.1 Evaluating Partitioning Analyzer

We now evaluate O²P on multiple benchmark datasets and workloads. Figure 1(a) shows the number of iterations in different variants of O²P for different tables in Star Schema Benchmark (SSB). We can see that O²Pp indeed improves over NV/HC on this realistic workload. O²Ppg and O²Ppgd are even better. Figure 1(b) shows the iterations in different variants of O²P for TPC-H dataset. For Lineitem table, O²Ppgd has just 42 iterations compared to 32,768 iterations in NV/HC. O²Ppgda and O²Ppgdm have the same number of iterations as O²Ppgd, hence we do not show them in the figure.

Next, we evaluate the actual running time of different O²P variants while varying the read-only workload. We vary the 100-query workload from OLTP style (1% tuple selectivity, 75-100% attribute selectivity) to OLAP style (10% tuple selectivity, 1-25% attribute selectivity) access patterns. We run this experiment over Lineitem (Figure 2(a)) and Customer tables (Figure 2(b)). We observe that on Lineitem O²Ppgd outperforms NV/HC by up to two orders of magnitude.

Now let us analyze the *quality* of partitioning produced by O²P. We define the quality of partitioning produced by an algorithm as the ratio of the expected query costs of optimal partitioning and the expected query costs of partitioning produced by the algorithm. The table below shows the quality and the number of iterations for optimal, NV, and O²P partitioning over mixed OLTP-OLAP workload. We can see that O²P significantly reduces the number of iterations, without losing much on partitioning quality.

Finally, we evaluate the scalability of O²P when increasing workload size. We vary the workload size from 1 to 10,000 queries consisting of equal number

	Customer			Lineitem		
	Optimal	Navathe	O2P	Optimal	Navathe	O2P
Quality	100%	99.29%	92.76%	100%	97.45%	95.80%
Iterations	100%	14.60%	2.28%	100%	2.42%	0.14%

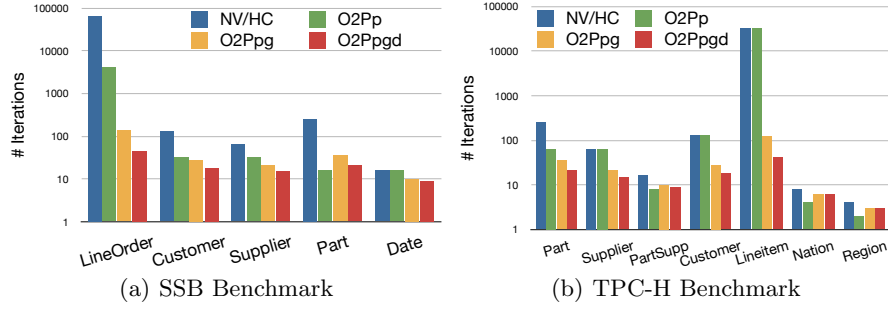


Fig. 1. Number of iterations in different algorithms over SSB and TPC-H benchmarks on different tables.

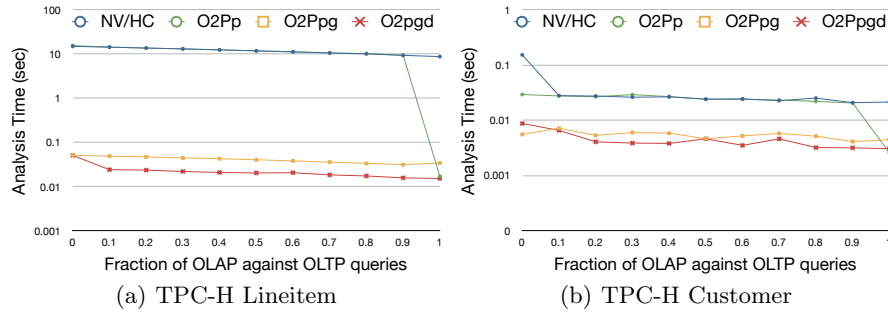


Fig. 2. Running times of different algorithms over changing workload type [100 queries each].

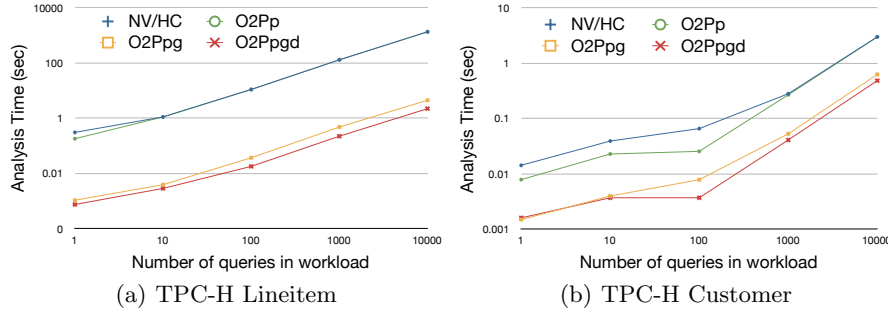


Fig. 3. Running time of different algorithms over varying workload size [with 50% OLAP, 50%OLTP queries].

of OLTP and OLAP-style queries. Figures 3(a) and 3(b) show the scalability of O^2P over TPC-H Lineitem and Customer tables respectively. We can see that all variants of O^2P algorithm scale linearly with the workload size. Hence, from now on we will only consider O^2Ppgd algorithm.

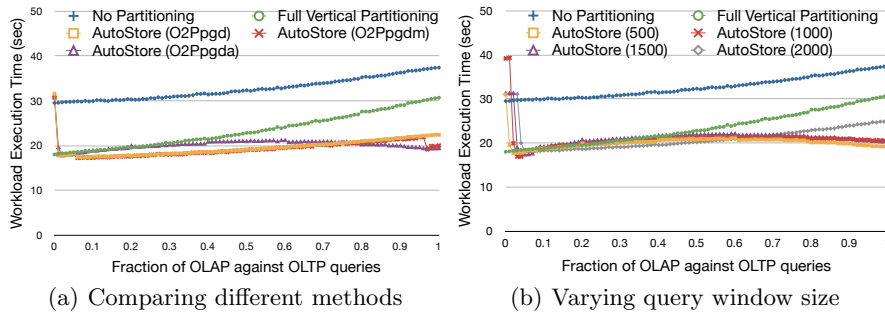


Fig. 4. Comparison of No Partitioning, Full Vertical Partitioning, and AUTOSTORE in main-memory implementation.

5.2 Evaluating Query Performance

Now we evaluate the query execution performance of AUTOSTORE in comparison with No and Full Vertical Partitioning. In this evaluation we use a main-memory implementation of AUTOSTORE in Java. In order to show how AUTOSTORE adapts vertical partitioning to the query workload, we use a universal relation de-normalized from a variant of the TPC-H schema [28]. Similar as in [28], we choose a vertical partition with part key, revenue, order quantity, lineitem price, week of year, month, supplier nation, category, brand, year, and day of week for our experiments. Further, since we consider equal size attributes only, we map all attributes to integer values, while preserving the same domain cardinality. We use a scale factor (SF) of 1.

Figure 4(a) shows the performance of No Partitioning, Full Vertical Partitioning, AUTOSTORE with O²Ppgd, AUTOSTORE with O²Ppgdm and AUTOSTORE with O²Ppgda. We vary the fraction of data accessed, i.e. both the attribute and tuple selectivity along the x-axis. We vary the OLTP/OLAP read access patterns as in Section 5.1, with a step size of 0.01%. From the figure we can see that AUTOSTORE automatically adapts to the changing workload, i.e. even though it starts with no-partitioning configuration, AUTOSTORE matches or improves full vertical partitioning performance. Therefore, from now on we consider only O²Ppgda. Figure 4(b) shows the performance of AUTOSTORE when varying query window size. From the figure we observe that larger query windows, e.g. query window of 2000 after 70% OLAP, become slower. This is because the partitioning analyzer has to now estimate the costs of more number of queries while analyzing partitioning schemes.

5.3 Evaluation over Real System

Modern database systems, e.g. PostgreSQL, have a very strong coupling between their query processors and data stores. This makes it almost impossible to replace the underlying data store without touching the entire software stack on top. This

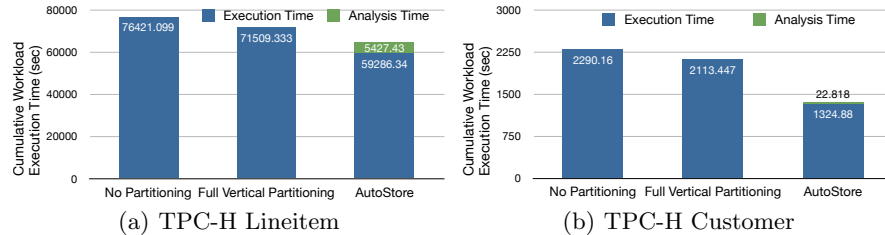


Fig. 5. Total running time of varying OLTP-OLAP workload on different tables in BerkeleyDB.

limitation led us to consider BerkeleyDB (Java Edition), which is quite flexible in terms of physical data organization, for prototyping.

In BerkeleyDB we store a key-value pair for each tuple. The key is composed of the partition ID and the tuple ID, and the value contains all attribute values in that partition. Since BerkeleyDB sorts the data on keys by default, we simply need to change the partition ID to change the partitioning scheme. Again, we vary the OLTP/OLAP read access patterns as in Section 5.1, with a step size of 0.01% over TPC-H dataset (total size 1GB). For different layouts, Figure 5(a) shows the total query execution times over TPC-H Lineitem table and Figure 5(b) shows the total query execution times over TPC-H Customer table. In general, AUTOSTORE outperforms the best layout in each of the tables. For instance, even though AUTOSTORE starts from a no-partitioning configuration, it improves over Full Vertical Partitioning by 36% in Customer table.

6 Related Work

Offline Horizontal and Vertical Partitioning. Horizontal partitioning is typically done based on values (range, hash, or list). A recent work proposed workload-based horizontal partitioning [12]. However, it is still offline. Vertical partitioning started with early approaches of heuristic based partitioning [16] of data files. The state-of-the-art work in vertical partitioning [23] develops the notion of attributes affinity, quantifying attribute co-occurrence in a given set of transactions. This work creates a clustered attribute affinity matrix in the first step and applies binary partitioning repetitively in the second step. A follow-up work [24] presents graphical algorithms to improve the complexity of their approach. Other works took the type of scan into account to analyze the disk accesses [10] and formulated an integer linear programming problem to arrive at the optimal partitioning for relational databases [10]. Next, researchers proposed transaction based vertical partitioning [9], arguing that since transactions have more semantic meaning than attributes, it makes more sense to partition attributes according to a set of transactions. However, all of these works considered data partitioning as a one-time *offline* process, in contrast to the online approach to data partitioning in AUTOSTORE. Recent works integrate parti-

tioning into physical database design tuning problem [3] with the objective of finding the configuration producing the minimum workload cost within the storage bound. This work first produces *interesting column groups* by applying a heuristic-based pruning of the set of all column groups. The column groups are then merged before all possible partitioning schemes are enumerated. These steps are not feasible in an online setting. HYRISE [15] analyzes the access patterns to partitions data. However, this approach is (1) still offline, i.e. it is not able to adapt to changing workloads, (2) restricted to main memory, whereas AUTOSTORE works for both disk and main memory DBMSs, and (3) is limited to vertical partitioning, whereas AUTOSTORE solves VPP and HPP equivalently.

Online Physical Tuning. Dynamic materialized views [31] materialize the frequently accessed rows dynamically. [6] proposes online physical tuning for indexes, without looking at the partitioning problem. Database Cracking [19] dynamically sorts the data in column stores based on incoming queries. However, these works still do not address the online partitioning problem.

7 Conclusion

In this paper, we revisited database partitioning with the objective of automatically fitting data to queries to an online query workload. We presented AUTOSTORE, an online self-tuning database store. AUTOSTORE monitors the workload and takes partitioning decisions automatically. We generalized VPP and HPP to the 1DPP. We presented the O^2P algorithm to effectively solve 1DPP. We performed an extensive evaluation of our algorithms over TPC-H data. We showed experimental results from a main-memory and a BerkeleyDB implementations of AUTOSTORE over mixed workloads. Our results show that O^2P is faster than earlier approaches by more than two orders of magnitude, and still produces good quality partitioning results. Additionally, our results show that over changing workloads AUTOSTORE outperforms existing stores.

Acknowledgements. Work partially supported by DFG, M2CI.

References

1. S. Agarwal et al. Database Tuning Advisor for Microsoft SQL Server 2005. In *VLDB*, 2004.
2. S. Agrawal, E. Chu, and V. Narasayya. Automatic Physical Design Tuning: Workload as a Sequence. In *SIGMOD*, 2006.
3. S. Agrawal et al. Integrating Vertical and Horizontal Partitioning into Automated Physical Database Design. In *SIGMOD*, 2004.
4. I. Alagiannis et al. An Automated, Yet Interactive and Portable DB Designer. In *SIGMOD*, 2010.
5. N. Bruno and S. Chaudhuri. Physical Design Refinement: The "Merge-Reduce" Approach. In *EDBT*, pages 386–404, 2006.
6. N. Bruno et al. An Online Approach to Physical Design Tuning. In *ICDE*, 2007.
7. N. Bruno and S. Chaudhuri. Constrained Physical Design Tuning. *PVLDB*, 2008.

8. S. Chaudhuri and V. Narasayya. An Efficient Cost-driven Index Selection Tool for Microsoft SQL Server. In *VLDB*, 1997.
9. W. W. Chu and I. T. Jeong. A Transaction-Based Approach to Vertical Partitioning for Relational Database Systems. *IEEE TSE*, 19(8):804–812, 1993.
10. D. W. Cornell and P. S. Yu. A Vertical Partitioning Algorithm for Relational Databases. In *ICDE*, 1987.
11. D. W. Cornell and P. S. Yu. An Effective Approach to Vertical Partitioning for Physical Design of Relational Databases. *IEEE TSE*, 16(2):248–258, 1990.
12. C. Curino et al. Schism: a Workload-Driven Approach to Database Replication and Partitioning. In *PVLDB*, 2010.
13. J. Dittrich and A. Jindal. Towards a One Size Fits All Database Architecture. In *CIDR*, 2011.
14. J.-P. Dittrich, P. M. Fischer, and D. Kossmann. AGILE: Adaptive Indexing for Context-aware Information Filters. In *SIGMOD*, 2005.
15. M. Grund et al. HYRISE - A Main Memory Hybrid Storage Engine. *PVLDB*, 2010.
16. M. Hammer et al. A Heuristic Approach to Attribute Partitioning. *ACM TODS*, 1979.
17. R. A. Hankins and J. M. Patel. Data Morphing: An Adaptive, Cache-Conscious Storage Technique. In *VLDB*, 2003.
18. J. A. Hoffer and D. G. Severance. The Use of Cluster Analysis in Physical Database Design. In *VLDB*, 1975.
19. S. Idreos et al. Database Cracking. In *CIDR*, 2007.
20. C. Jermaine, E. Omiecinski, and W. G. Yee. The partitioned exponential file for database storage management. *The VLDB Journal*, 16:417–437, 10 2007.
21. A. Jindal. The Mimicking Octopus: Towards a one-size-fits-all Database Architecture. In *VLDB PhD Workshop*, 2010.
22. H. Kimura et al. CORADD: Correlation Aware Database Designer for Materialized Views and Indexes. In *VLDB*, 2010.
23. S. Navathe, et al.. Vertical Partitioning Algorithms for Database Design. *ACM TODS*, 1984.
24. S. Navathe and M. Ra. Vertical Partitioning for Database Design: A Graphical Algorithm. In *SIGMOD*, 1989.
25. P. E. O’Neil, et al. The Log-Structured Merge-Tree (LSM-Tree). *Acta Inf.*, 1996.
26. O. Ozmen, K. Salem, J. Schindler, and S. Daniel. Workload-Aware Storage Layout for Database Systems. In *SIGMOD*, 2010.
27. S. Papadomanolakis and A. Ailamaki. AutoPart: Automating Schema Design for Large Scientific Databases Using Data Partitioning. In *SSDBM*, 2004.
28. V. Raman et al. Constant-Time Query Processing. In *ICDE*, 2008.
29. D. Sacca and G. Wiederhold. Database Partitioning in a Cluster of Processors. *ACM TODS*, 10(1):29–56, 1985.
30. K. Schnaitter et al. COLT: Continuous On-Line Database Tuning. In *SIGMOD*, 2006.
31. J. Zhou et al. Dynamic Materialized Views. In *ICDE*, 2007.
32. D. C. Zilio et al. DB2 Design Advisor: Integrated Automatic Physical Database Design. In *VLDB*, 2004.