

An Adaptive Online System for Efficient Processing of Hierarchical Data *

Athanasia Asiki

Dimitrios Tsoumakos

Nectarios Koziris

Computing Systems Laboratory
School of Electrical and Computer Engineering
National Technical University of Athens
{nasia, dtsouma, nkoziris}@cslab.ece.ntua.gr

ABSTRACT

Concept hierarchies greatly help in the organization and reuse of information and are widely used in a variety of information systems applications. In this paper, we describe a method for efficiently storing and querying data organized into concept hierarchies and dispersed over a DHT. In our method, peers individually decide on the level of indexing according to the granularity of the incoming queries. Roll-up and drill-down operations are performed on a per-node basis in order to minimize the required bandwidth for answering queries on variable aggregation levels. We motivate our approach by applying it on a large-scale Grid system: Specifically, we plan to apply our fully decentralized scheme that creates, queries and updates large volumes of hierarchical data on-line and replace the traditional centralized and strictly indexed information systems. Our extensive experimental results support this argument on many diverse configurations: Our system proves very efficient in skewed workloads, both over single and multiple hierarchy levels at the same time. It adapts to sudden changes in popularity and effectively stores and updates large amounts of data at very low cost.

Keywords

Distributed Hash Table, Concept Hierarchies, Adaptive Indexing, Grid Information System

1. INTRODUCTION

A concept hierarchy (or *taxonomy*) defines a sequence of mappings from more general to lower-level concepts. For example, Figure 1 shows a simple hierarchy for the Virtual Organization concept, where $VO < Category < Region < Site$ and one for time where a *partial* order is defined. Concept hierarchies are important because they allow the structuring of information into

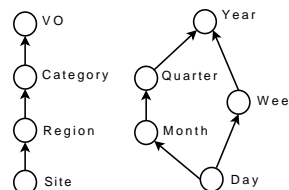


Figure 1: A concept hierarchy for the VO and Time (lattice) dimensions.

categories, thus enabling its search and reuse. Specifically, users may view data at different levels of a dimension hierarchy: With the *roll-up* operation we climb up to a more summarized level of the hierarchy, while a *drill-down* defines the opposite operation (i.e., navigating to lower levels of the hierarchy with increased detail). The drilling paths are usually defined by the hierarchies within the dimensions. The mappings of a concept hierarchy are usually provided by application or domain experts.

Works in the field of data-warehousing (e.g., [10, 18], etc) utilize hierarchies across the dimensions of a data cube but these present strictly centralized solutions. In the area of distributed computing, while there has been considerable work in sharing simple relational data using both structured and unstructured overlays (e.g., [13, 14, 17]), no special consideration has been given to data supporting hierarchies. We investigate the problem of indexing and querying such data in a way that preserves the semantics of the hierarchies and is efficient in retrieving the requested values in a fully distributed environment.

To motivate our approach, we describe how it can be applied in order to function as a distributed and efficiently operational grid *information system*. Grid computing allows for coordinated resource sharing and problem solving in dynamic virtual organizations (VOs). A VO is a group of users from multiple institutions who collaborate to achieve a specific goal. The goal of grid computing is to provide a network of systems which, acting like a single supercomputer, offers resources that are easily accessible. In order for jobs to be adequately served by the most appropriate resources, the information system stores all needed information about the characteristics of the available resources over time.

There exist a number of systems that accomplish the tasks of an information system (e.g., [3, 5], etc). Nevertheless, they either feature a central information repository or a hierarchy of aggregation sites that introduce both scalability (single points of failure) and performance (processing burden on a single site) issues. Work in the area of distributed databases offers a variety of systems which can be used to disseminate and query this information. Nevertheless, these schemes cannot be used to maintain the semantics of the hierarchy and efficiently retrieve views of the data at different

*©ACM, (2009). This is the author's version of the work. It is posted here by permission of ACM for your personal use. Not for redistribution. The definitive version was published in Proceedings of HPDC'2009 <http://doi.acm.org/10.1145/nnnnnn.nnnnnn>

granularities. This is very important for applications such as a large scale information system, as queries naturally target different levels of detail: Historic queries usually require grouping by the highest hierarchy levels (e.g., group-by VO or group-by Year), whereas online queries are naturally directed towards more detailed levels (e.g., group-by Site or by Day).

Let us assume that the system's database contains a location dimension that relates to the VO location information (see Figure 1). Monitoring information is described by attributes (*facts*) such as the number of running jobs, number of waiting jobs, available storage space, total storage space, etc. Common accounting queries could be "Give me the average CPU time" or "Give me the minimum available space in Gbytes", presented to the user grouped by a VO value.

In this paper, we present a system that efficiently stores, queries and updates data organized in concept hierarchies. We choose the DHT as a reliable substrate over which we store, index and query our data, thus eliminating all possible bottlenecks created by the hierarchical or centralized approaches mentioned before. Data producers individually insert data to the distributed information system. Queries are still answered, while incremental updates are efficiently processed. Our solution takes the query granularities into account, adjusting the indexing structure to favor performance. Second, we intend to provide a system that will preserve all hierarchy-specific information. In our technique, a tree-like data structure is used to store data and maintain indices to related keys, enabling us to respond to more complex, hierarchy-based queries such as: "Which sites correspond to VO 'Biomed' " or "What category does region 'AsiaPacific' belong to ". We can summarize the contributions of this paper in the following points:

- We present a complete storage, indexing and query processing system for hierarchical data. This system has many desirable properties: It adapts the granularity of its indexing according to incoming requests; Performs efficient and online incremental updates; Maintains data in a fault-tolerant and fully distributed environment.
- We motivate the usefulness of this scheme by customizing it to serve as a high-performance information system. We show how our method outperforms traditional approaches by eliminating offline processing and other performance bottlenecks.
- We present a thorough performance analysis in order to identify the behavior of our scheme under a large range of work and data loads.

The rest of this paper is organized as follows: The next Section summarizes related work both in exploiting hierarchies as well as existing information systems. Section 3 defines the problem and presents our solution in detail, while Section 4 refers to the case study of the information system. Section 5 describes our experimental setup and the collected results, while we conclude our work in Section 6.

2. RELATED WORK

There has been significant work in the area of databases over P2P networks. PIER [13] proposes a distributed architecture for relational databases supporting operators such as join and aggregation of stored tuples. A DHT-based overlay is used for query routing. The Chatty Web [7] considers P2P systems that share (semi)-structured information but deals with the degradation, in terms of syntax and semantics, of a query propagated along a network path.

In [20], the authors propose optimization techniques for query reformulation in P2P database systems.

In GrouPeer [14], SPJ queries are sent over an unstructured overlay in order to discover peers with similar schemas. Peers are gradually clustered according to their schema similarity. PeerDB [17] also features relational data sharing without schema knowledge. Query matching and rewriting is based on keywords provided by the users. GridVine [8], and pSearch [19] are based on structured P2P overlays. GridVine hashes and indexes RDF data and schemas, and pSearch represents documents as well as queries as semantic vectors. All these approaches offer significant and efficient solutions to the problem of sharing structured and heterogeneous data over P2P networks. Nevertheless, they do not deal with the special case of hierarchies over multidimensional datasets.

An interesting method for representing hierarchical data is presented in [15]. The method is applied on unstructured networks containing XML documents in order to favor the routing of path queries. Each XML document is represented by an unordered label tree and bloom filters are used to summarize it.

Several indexing schemes have been presented for storing data cubes (e.g., [16, 21]). However, only few support both aggregate queries and hierarchies. In [18], hierarchies are exploited to enable faster computation of the possible views and a more compact representation of the data cube. The *Hierarchical Dwarf* contains views of the data cube corresponding to a combination of the hierarchy levels. The other approach is the DC-Tree [10]. In this work, the attributes of a dimension are partially ordered with respect to the valid hierarchy schema for each dimension. The DC-tree stores one concept hierarchy per dimension and assigns an ID to every attribute value of a data record that is inserted. These approaches are very efficient in answering both point and aggregate queries over various data granularities but do so in a strictly centralized and controlled environment.

There exist multiple systems and architectures proposed to implement the information system component. The most common is the *Globus Monitoring and Discovery Service (MDS)* [3]. A *Grid Index Information Service (GIIS)* provides an aggregate directory of lower level data stored at multiple *Grid Resource Information Services (GRISs)*. The hierarchical structure that can be composed between GIISs enables complete information retrieval by querying the top level GIIS. However, MDS has shown not to be a solution for large-scale production because it does not scale: Multiple client requests quickly lead to an overload of the top level GIIS [22]. Another schema used especially for accounting and publication of user-level information is the *Relational Grid Monitoring and Discovery Service (R-GMA)* [5], which supports complex type of queries allowed by relational databases. R-GMA presents information as a single virtual database containing a set of virtual tables, nevertheless the bulk of data need be transferred offline to a centralized database after a period of time, with all the performance drawbacks that this entails. The major drawback in all these methods is the fact that none of these architectures scale as the number of data collectors increase [22]. Moreover, they all assume an offline (or periodic at best) data migration phase to a more central location where global information can be available. In contrast, we propose a completely decentralized system where all data are continuously available and indexed at the requested granularities for fast retrieval.

3. AN ADAPTIVE INDEXING SCHEME TO SUPPORT CONCEPT HIERARCHIES

In this paper, we describe a system for processing bulk hierar-

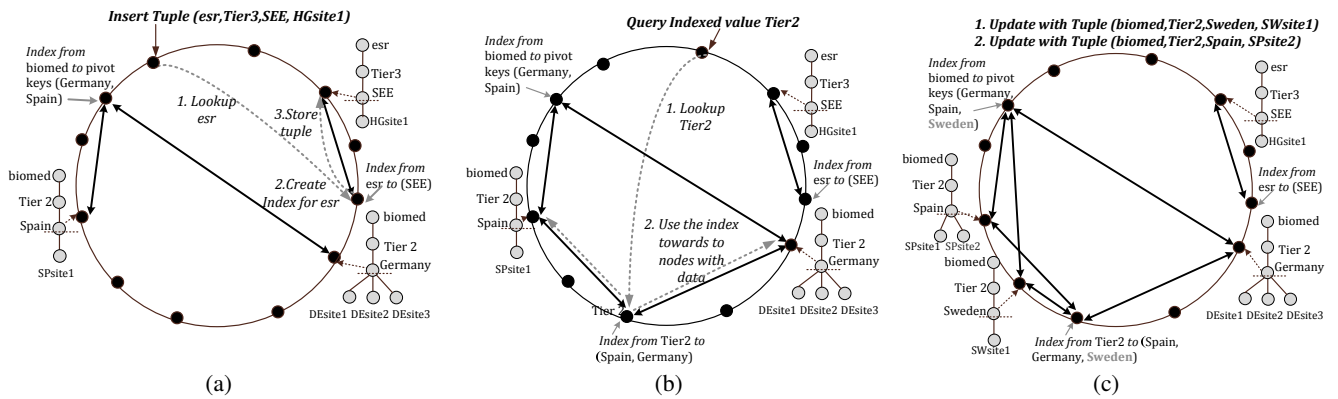


Figure 2: (a) Insertion of a new tuple with its root key not already stored in the overlay (b) Lookup using soft-state indices for a value belonging to a non-pivot level (c) Updates for an already existing and a non-existing pivot key

chical data in a DHT-based overlay. Our goal is to enable efficient querying while preserving the hierarchy semantics. In addition to the fact that the DHT substrate can transparently handle node churn, replication, reliable distributed storage, etc, our technique offers adaptive indexing according to the granularity of the incoming queries and online updating with low cost and no downtime.

3.1 Notation

Let the data items stored in the system be in the form of *tuples* containing values for all levels l_i of a concept hierarchy with L levels. They also contain a numerical fact of interest (e.g., CPU Time, Available Memory, etc) or the location of actual data. We call the uppermost level of the hierarchy (l_0) *root level* and its value *root key*. We define that $l_a < l_b$, where $(a, b \in [0, L - 1])$, if and only if l_a is higher in the concept hierarchy (namely closer to l_0) than l_b . The values of the hierarchy levels are organized in tree structures, one per root key. Without loss of generality, we assume that each value of l_i has at most one parent in l_{i-1} . During the insertion of a tuple, a level of its hierarchy is chosen and its hashed value serves as its *key* in the underlying DHT overlay. We refer to this level as *pivot level* and to its value as *pivot key*. Finally, the highest and the lowest pivot levels of the hierarchy for a specific root key are called *MinPivotLevel* and *MaxPivotLevel* respectively.

3.2 Data Insertion

The key for each tuple is the result of a hash function applied on the value of the selected pivot level. Tuples are assigned to the node with ID numerically closest to the generated keys, according to the standard DHT operations.

In our system, both initial insertions as well as incremental updates are handled in a unified manner. We introduce a completely distributed catalogue containing all the root keys and their corresponding pivot keys. Each root key is stored in the node responsible for it along with the list of pivot keys that have been already inserted. The root key is also aware of the *MaxPivotLevel* used during tuple insertion containing its value.

The procedure followed during the tuple insertion is as follows: The root key for this tuple is generated and a lookup for this key takes place in the DHT overlay. If the root key exists, the tuple ends up in the node responsible for it. We consider an *insertion* to be the procedure followed in case that the root key does not already exist in the overlay. Otherwise, the *update* procedure is followed (see Section 3.5).

In case of an insertion, the tuple or the group of tuples with the same root key arrive at the node responsible for it. The node selects a pivot level (either a random or a predefined one) and the pivot

key(s) of the tuple(s) is (are) calculated for the appropriate pivot level. Each new pivot key is added to the list of pivot keys. Finally, each tuple is stored in the node with the ID closest to its pivot key.

Each peer organizes the tuples in trees that preserve their hierarchical nature. As a consequence, each distinct value of the pivot level corresponds to a tree that reveals part of the hierarchy. When a new tuple arrives at the node responsible for it, the node searches its keys. If no tuples for this pivot key have been stored, a new tree with a single branch is created. In the opposite case, a new branch is added below the value of the pivot level with the new values of the remaining levels.

An example of insertion is shown in Figure 2(a). A graphic convention is that solid lines represent existing indices, while dotted lines correspond to logical steps followed during the described procedure. Let us assume that tuples follow the VO hierarchy depicted in Figure 1 with l_2 (Region) as the globally defined pivot level for initial insertions. A tuple with root key 'esr' is inserted in the overlay. Since the specific root key does not already exist, a new *index* is created in the corresponding node. Afterwards, Region is selected as pivot level and the tuple is forwarded to the node responsible for this key, which creates a new tree with only one branch for this tuple.

3.3 Data Lookup and Soft-state Indices

Queries concerning the pivot level are exact match queries and can be answered by the DHT lookup operation. Queries for any other level cannot be resolved unless flooded across the DHT. Towards the exploitation of the knowledge acquired by flooded queries, we introduce *soft-state bidirectional indices* to our scheme. When a node answers a flooded query, it checks whether a roll-up or drill-down is necessary. If this is not the case, the query initiator starts the procedure of creating an index, as soon as it receives the complete answer. It inserts the result of hashing the requested value in the DHT along with IDs of nodes having the actual tuples. The tuple holders also mark the specific value as *indexed*. The next time that a query for this key is initiated, the lookup operation locates the node holding the index, finds out the IDs of nodes with relevant tuples and retrieves them.

The created indices are soft-state, in order to minimize the redundant information. This means that they expire after a predefined period of time (Time-to-Live or *TTL*). Each time that an existing index is used, its *TTL* is renewed. This constraint ensures that changes in the system (e.g., data location, node departures, etc) will not result in stale indices, affecting the validity of the lookup mechanism. While memory becomes a cheaper commodity by the

day, the plain size of data discourages an “infinite” memory allocation for indices. Therefore, after the number of indices has reached a limit I_{max} , the creation of a new index results in the deletion of the oldest one. Calibrating I_{max} for performance without increasing it uncontrollably entails knowledge of our data (e.g., how skewed each hierarchy is). Overall, the system tends to preserve the most “useful” indices, namely the ones directed towards the most frequently queried data items.

The nodes with actual tuples of the indexed value need to know the existence of an index, in order to erase it after a re-indexing operation. The bidirectionality of the indices is introduced only to ensure data consistency, despite of them being soft-state. During re-indexing operations, the locations of stored tuples change and indices correlated to these tuples need either to be updated or erased, preventing the existence of stale indices. Our choice is to erase them, so as to avoid increasing the complexity of the system. Detailed information for an existing index is not essential for the node, where the tuples are stored. A simple mark for each indexed value is adequate in order to erase its index, if it is needed. In this case, some redundant operations for erasing expired indices may occur. If there are no memory restrictions and local processing is preferable to bandwidth consumption, indexed values can be marked with a time-stamp. Every lookup for an indexed value renews the TTL in both sides of the index and only valid indices are erased during re-indexing operations.

In the example of Figure 2(b), a query for ‘SEE’ is resolved directly by the lookup operation of the DHT protocol. Lookups for values of the root level are processed utilizing the indices created during insertions. Nevertheless, any query for any level other than the pivot level or the root level ends up with no results. For example, a query for date items described by ‘Tier2’ does not contact the two nodes with the corresponding trees before the creation of the index. The next step is the flooding of the query and the nodes storing the keys of ‘Spain’ and of ‘Germany’ are reached. The query initiator, which is now aware of the existing pivot keys, creates an index and these pivot keys are stored to the node responsible for the value ‘Tier2’. This node now has an index pointing to the node ‘Spain’ and another to node ‘Germany’. In the future, queries for ‘Tier2’ will be answered without flooding, utilizing the created soft-state indices. As shown analytically in Figure 2(b), a subsequent query for the Category ‘Tier2’ reaches the node responsible for this key, which in turn forwards the query to all relative nodes directly. The nodes storing the trees with the queried value return only the relevant tuple(s).

3.4 Re-indexing Operation

Our goal is the described system to dynamically adapt on a per node basis to online queries, so as to increase the ratio of the non-flooded queries. In order to achieve this goal, we introduce two re-indexing operations regarding the selection of pivot level: *roll-up* towards more general levels of the concept hierarchy and *drill-down* to levels lower than the pivot level.

The idea behind individual re-indexing of stored tuples is based on the fact that each node has a global view of the queries regarding each level $\ell_i < pivotlevel$, but only a partial view of the queries for each level $\ell_i > pivotlevel$. Therefore, it has sufficient information to decide if a drill-down will favor the increase of the exact-match queries for its values. On the other hand, a node has to cooperate with other peers that store a value of a level $\ell_i < pivotlevel$ in order to decide if this level is more appropriate.

The re-indexing of the data tuples (through a choice of a different pivot level) is performed on a per-tree basis, requiring no global coordination. Each node collects information using the incoming

queries and finds out if the pivot level of a tree remains its most popular level. Otherwise, the node proceeds with the re-indexing of the tuples of this tree. The popularity of the levels of a tree is estimated based on their average rates of incoming queries (hence InQ). A node maintains one record per tree with these rates during a restricted time-frame W . This parameter should be properly selected to perceive variations of query distributions and, at the same time, stay immune to instant surges in load.

In more detail, the mechanism works as follows: A node may check if a re-indexing is required based on the objective to achieve. The implemented strategy implies that a node decides whether a roll-up or drill-down is required when it answers a flooded query or when a number of queries for indexed values have been received. While the main objective is the increase of the queries answered without flooding, this strategy targets to the increase of exact match queries as well.

The number of queries for indexed values triggering a node to examine a possible re-index may vary and has an impact to the adaptiveness of the system. A small value indicates that the potential of re-indexing is examined more often and thus more re-indexing operations may take place. Nevertheless, if a decision has erroneously been made, it can be easily corrected. However, during re-indexing operations, existing indices are deleted and this may have a negative impact on the system. In the opposite case, the system tends to depend more on the effectiveness of the indices. We have observed that re-indexing operations are necessary, when popular values belong to levels with several distinct values. Indices perform better for higher levels with less values, since there is a high probability for repeated utilization of an index.

A node decides if a re-indexing operation will favor the increase of non-flooded queries based on the ‘popularity’ of each level according to the procedure described in basic steps in Algorithm 1. A *thr* parameter is used to indicate if a re-indexing operation is required. The following criterion defines if a re-indexing to a level ℓ_q is allowed:

$$InQ_{\ell_q} > thr \times \sum_{i=0}^{i=L-1} InQ_{\ell_i} \text{ where } \ell_q \neq pivotlevel, thr \in [0, 1]$$

In the described algorithm, two possible cases are considered to indicate the necessity of a re-indexing operation: **The queried level ℓ_q lies lower in the hierarchy than the pivot level of the tree ($\ell_q > pivotlevel$).** Only one tree stores the values of a level below the pivot level. Therefore, the specific node is aware of the exact popularity of these values and feels ‘confident’ to decide if a drill-down is needed. If the most popular level ℓ_{pop} of the tree lies below the pivot level and the defined criterion is valid for its InQ , then a drill-down to this level is decided. After the decision for drill-down is made, the node finds all the distinct values of the new pivot level and hashes them one by one, sending the new groups of tuples to the corresponding nodes. The already gathered statistic information is sent along with one randomly selected group, in order to maintain information about the query distribution for the values contained in the drilled-down tree within W . Any existing indices for any value of this tree are removed. If a drill-down is not needed, the node includes in its answer to the initiator the fact that the queried level is lower than the pivot level, hence it can carry on with the creation of the soft-state index and expedite the process.

The queried level ℓ_q lies higher than the pivot level of the tree ($\ell_q < pivotlevel$). In this case, there are more than one trees with this value needed to participate in a possible roll-up to this level. Otherwise, lookups for this value will not return complete results. If the threshold criterion is satisfied for the ℓ_q , then the node is positive to the potential of adopting this level as pivot level for this

tree. This step is indicative of an imbalance and the query initiator is informed about this. The query initiator decides for a re-indexing operation according to the procedure described in Algorithm 2. If the query initiator is aware of at least one node willing to roll-up to this level, it starts a procedure to confirm the local intuition by using statistic information provided by all the nodes having answered the query. After receiving the tuples containing the number of InQ per level, it calculates the total value of InQ per level.

The calculation of the total rate of InQ per level is not straightforward. Queries concerning an ℓ_i for any $i \geq \text{pivotlevel}$ end up only in one node and are thus counted once for statistic purposes. The same property is not valid for queries requiring values of higher levels than the pivot level. These queries reach more than one node and are counted in all of them. During the gathering of statistic information for roll-up decisions, the problem of multiple counting of such queries in the calculation of the rate for each level needs to be solved. The complexity in the calculation of the overall rate of InQ increases since more than one pivot levels may exist for the involved trees in the re-indexing procedure. For example, let us assume that the state of trees with the 'biomed' as their root key is the one shown in Figure 3(b). In this case, the value of InQ for 'Tier2' is sent twice by the nodes with the trees of 'Tier2'. To avoid this situation, a path containing the values for all levels in $[0, \text{pivotlevel}]$ is sent along the statistic information to the querying node, so as to make the correct decision. Through this procedure, the querying node is also informed for the *MinPivotLevel* and *MaxPivotLevel* of all existing trees containing the queried value (hence *NotPivotKey*).

If the InQ of ℓ_q is more than *thr* of the total number of InQ, then the initiator messages the involved nodes to roll-up the corresponding trees to this level by re-inserting their tuples. If the re-indexing criterion for ℓ_q is not fulfilled and since statistic information has been collected, the querying node examines if a drill-down to a level $\ell_i \geq \text{MaxPivotLevel}$ (the equality is for the case that all the involved trees do not have the same pivot level) is dictated by the collected statistics. Our intention is to take advantage of the fact that the querying node has now a more global view of the InQ per level. It is possible to find a level $\ell_i \geq \text{MaxPivotLevel}$ to be the most popular but this tendency not to appear in the partial views of the involved nodes. In this case, the query initiator informs the involved nodes that a drill-down to this level is needed. We call this procedure *Group-Drill-down*, since more than one nodes participate in the drill-down. All the trees with the queried value in ℓ_q drill-down to the new pivot level. If the new pivot level is equal to the *MaxPivotLevel*, the trees already in the *MaxPivotLevel* do not perform any action. If a re-indexing operation is not needed, no action is taken other than the creation of a soft-state index for this value.

Lock mechanisms are activated during the time that a re-indexing decision is being made. The purpose of this locking is to avoid examining the same re-indexing possibility multiple times for concurrent lookups on specific trees. Locks are revoked after the completion of an ongoing procedure or after a short period of time. The steps described in Algorithms 1 and 2 are performed, only if the corresponding locks are inactive. Otherwise the described procedures are not performed and only the query is answered.

Examples of the described re-indexing operations are applied in the trees of Figure 3 with root value 'biomed'. The two trees of Figure 3(a) are stored in different nodes of the DHT overlay and are considered the initial state before any re-indexing operation. We suppose that a query for 'Spain' triggers a drill-down operation. The result is shown in Figure 3(b). On the other hand, a query for 'biomed' may result in a roll-up to the root level depicted in Figure 3(c) or a *Group-Drill-down* to the *Region* level depending on the

Algorithm 1 Decision Algorithm in the node answering a query

```

pivotlevel: current pivot level
 $\ell_q$ : the queried level of the flooded or indexed value
NotPivotKey: the flooded or indexed value
 $InQ_{tot}$ : rate of incoming queries for the tree with NotPivotKey
 $InQ_{ini}$ : initial minimum rate to allow re-index operations
 $InQ_{l_{pop}}$ : rate of incoming queries for the most popular level
action: the decided action
 $\ell_{pop} \leftarrow \text{FindMostPopularLevel}$ 
if  $\ell_q > \text{pivotlevel}$  then
  if  $(InQ_{tot} > InQ_{ini})$  AND  $(\ell_{pop} > \text{pivotlevel})$  AND
 $(InQ_{\ell_{pop}} > thr \times InQ_{tot})$  then
    Drill-down to  $l_{pop}$ 
    action  $\leftarrow$  NoAction
  else if NotPivotKey is NOT indexed then
    action  $\leftarrow$  CreateIndex
  else
    action  $\leftarrow$  NoAction
  end if
else if  $\ell_q < \text{pivotlevel}$  then
  if  $(InQ_{tot} > InQ_{ini})$  AND  $(\ell_q = \ell_{pop})$  AND  $(InQ_{\ell_q} > thr \times$ 
 $InQ_{tot})$  then
    action  $\leftarrow$  PositiveToRollup
  else if NotPivotKey is NOT indexed then
    action  $\leftarrow$  CreateIndex
  else
    action  $\leftarrow$  NoAction
  end if
end if

```

total InQ per level. The *Group-Drill-down* results to the trees of Figure 3(d) and differs from the simple drill-down, since all the trees drill-down to ℓ_3 .

3.5 Updates

An *update* is the procedure followed during the insertion of a tuple, when its root key already exists in the overlay. The update procedure comprises of two consecutive phases: the insertion of the tuple and the updating of any existing indices for the values of the tuple. The insertion phase presents minor differences compared to the insert procedure.

The updates of the existing datasets is a more complicated procedure. During the insertion of new tuples, it is critical to select the correct pivot level so as to ensure the correctness of the lookup operations. The selection of the pivot level is not simple, since the pivot levels of the stored trees of a specific root key may vary due to performed re-indexing operations.

The following assumptions are made:

- If a new tuple contains a pivot key, then this key should be used during insertion. Otherwise, lookups for this value will return only the tuples, that existed before this operation.
- Even if none of the values belonging to the specific tuple have been used as pivot keys, they may have already been stored in the network. The selection of such a value as pivot key would result in the discovery of the new tuple only in a later search. Therefore, we consider that the pivot level be equal to the *MaxPivotLevel* in this case. Re-indexing operations would take over to find the most appropriate pivot level of this tuple.

The difficulty of updates increases because the information about the stored pivot values and the *MaxPivotLevel* is distributed over

Algorithm 2 Decision Algorithm in the querying node

l_q : the queried level of the flooded or indexed value
NotPivotKey: the flooded or indexed value
action: the required action by involved nodes {*action* = *PositiveToRollup* if at least one node is possitive to roll-up}
if *action* = *PositiveToRollup* **then**
 Gather statistic information
 Calculate InQ for each level
 $l_{pop} \leftarrow FindMostPopularLevel$
 $MaxPivotLevel \leftarrow FindMaxPivotLevel$
 if ($l_q = l_{pop}$) AND ($InQ_{l_{pop}} > thr \times InQ_{tot}$) **then**
 Roll-up to l_{pop}
 else if ($l_{pop} \geq MaxPivotLevel$) AND ($InQ_{l_{pop}} > thr \times InQ_{tot}$) **then**
 Group-Drill-down to l_{pop}
 else if *NotPivotKey* is NOT indexed **then**
 Create Index for *NotPivotKey*
 end if
else if *action* = *CreateIndex* **then**
 Create Index for *NotPivotKey*
end if

the overlay. We have implemented a distributed catalogue by creating an index among the node responsible for the root key and the nodes with the pivot keys, namely a record with the root key and the corresponding pivot keys. This enhancement allows online updates with the system continuing to efficiently serve requests of the users.

During a new tuple insertion, a lookup operation for the root key is performed. The responsible node is contacted and it finds out if any value of the tuple corresponds to a pivot key. In this case, the tuple is stored to the responsible node for the pivot key and its new values below the pivot level are added as a new branch to the existing tree. In the opposite case, the hashed value of the $MaxPivotLevel$ is considered as the pivot key of this tuple during its insertion in the overlay. The existence of trees with equal values above the pivot level is not excluded by this assumption and neither is the existence of corresponding indices. These indices should be updated so as indexed lookups to return complete answers. The node storing the tuple initiates lookups for each l_i , where $0 < l_i < pivotlevel$ and the corresponding indices are informed about the new tuple.

Examples for the possible cases during an update are depicted in Figure 2(c). The node holding the index for the root key '*biomed*' concludes that the none value of the tuple '*(biomed, Tier2, Sweden, SWsite1)*' corresponds to an existing pivot key. Moreover, it is aware that the $MaxPivotLevel$ for its trees is the *Region* and thus the tuple is inserted with '*Sweden*' as its pivot key. The new pivot key is also added in the list of pivot keys for this root key. However, the value '*Tier2*' is already indexed. During lookup for the value '*Tier2*' according the update procedure, the responsible node is discovered and a new index among this node and the node with actual data is created. During the update for tuple '*(biomed, Tier2, Spain, SPsite2)*', the node with '*biomed*' index proceeds in the insertion of tuple with '*Spain*' as pivot key, resulting to the creation of a new branch in the existing tree. The indexed value '*Tier2*' is not affected and no further action is taken.

4. CASE STUDY: GRID INFORMATION SERVICES

A motivating scenario for the usefulness of the proposed system

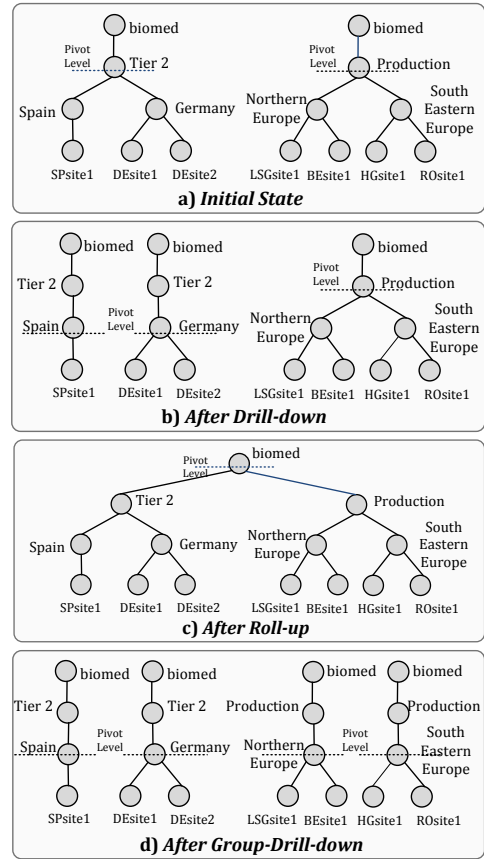


Figure 3: Examples of drill-down and roll-up operations

can be found in the collection of information produced by information services in Grid environments. Grid computing resources and services advertise a large amount of data, which are used by multiple people across multiple administrative domains. This information is not intended only for event handling, as in traditional monitoring systems for networks and cluster computing. The produced information by various mechanisms such as cluster monitors (Ganglia [2], Hawkey [4]), services (GRAM, RLS [6]), queuing systems, etc, is organized and provided to various applications, such as accounting systems, schedulers, portals, etc. For this reason, the key to the design of Grid Information Services is to identify the information that is required and to determine how to best make this information available.

In recent Grid Monitoring Architectures (MDS, R-GMA), data concerning the state of the infrastructure are collected by a combination of various monitoring systems on a resource base and organized by *information producers (or providers)*. In the previous generation of monitoring architectures, the information producers were organized in a hierarchical structure and published their data, which finally were collected and stored in a central LDAP-based database. In more modern approaches, the information producers are known to the system by subscribing themselves to an *Index service* (MDS) or a *Registry* (R-GMA). Information consumers ask this structure for the location of the producers and contact all of them in order to acquire the needed data. Moreover, various aggregator services exist that collect information (via subscription, polling or execution) from information producers using a common configuration mechanism to specify the type of data and the collected information. For example, VOs maintain such services to

maintain VO-wide resource information by collecting data from the Information servers running at many sites. Due to the large volume of data and their usefulness to various services, we strongly believe that a solution for efficient storage and indexing of the produced information adapting to the requests of users or services contributes to the operability and the performance of a Grid Infrastructure.

Our system is a complete solution for the organization and storage of such information. The proposed scheme can be integrated in a Grid environment providing a fully decentralized solution. In this architecture, the nodes hosting the information producers and used for information-related purposes are organized in a DHT-based overlay, as shown in Figure 4. The routing of messages among these resources is performed according to the DHT protocol and no centralized structures are required. The P2P overlay introduces a scalable solution, while data and query load balancing is achieved. A concept hierarchy is defined for the various levels of aggregation that characterize the produced data and the existing needs. Each numerical fact is described by the corresponding concept hierarchy. No off-line collection and processing of data is required, since online updates are supported. The re-indexing mechanism enables the summarization of data according to the incoming queries. The Registry or Index service with the locations of service instances is implemented in a distributed manner. Data are distributed among the nodes of the DHT and queried according to our method. While in existing approaches a consumer queries all the information producers, in our system only the nodes responsible for data according to the aggregation level are contacted reducing the communication and processing cost.

An example for the usage of the described integration can be considered for the service providing information for accounting purposes of the EGEE Accounting Portal [1]. This service uses APEL [11], which is a log processing application to filter data produced in each site. Afterwards, R-GMA producers collect data from sites and streams them to a centralized database. The central database collects millions of records per grid job and stores them offline. Due to the enormous volume of data, only summarized views of the various metrics such as Number of jobs, Normalized CPU usage, SumCPU, CPU efficiency, etc, are computed offline and become available to the users. Needless to say, the right balance between quantity and timeliness of information, on the one hand and associated costs on the other should be considered for this centralized solution. A concept hierarchy describing the measurements needed for the above metrics is the VO concept hierarchy. We selected this hierarchy, since this application allows queries with group by per VO, per Site, per Category and per Region (or per country). In our system, the offline database is distributed across the DHT overlay. A query of a user starts from any node and ends up in the node(s) with relative data, where it is being processed. Our indexing scheme adapts the summarization level of the produced records from more general ones (level VO) to more detailed ones (level Site) according the incoming queries.

5. EXPERIMENTAL RESULTS

5.1 Simulation Setup

We now present a comprehensive simulation-based evaluation of our scheme. Our performance results are based on a heavily modified version of the FreePastry simulator [12], although any DHT implementation could be used as a substrate. We assume a network size of 256 nodes, all of which are randomly chosen to initiate queries.

In our simulations, we use synthetically generated data. Our data is a tree with each value having at most one parent. Each distinct

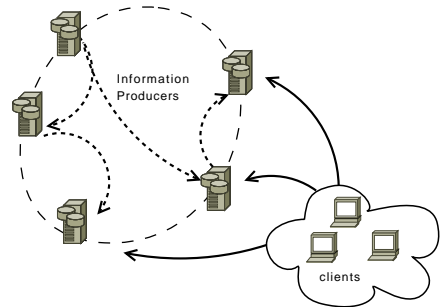


Figure 4: The proposed architecture for the Information System

value of ℓ_i has a constant number of children in ℓ_{i+1} . By default, our data comprise of 100k tuples, organized in a 4-level hierarchy (see Figure 1(a)) with one numerical fact (e.g., CPU_time). The number of distinct values per level are $|\ell_0 = 100|$, $|\ell_1 = 1000|$, $|\ell_2 = 10000|$ and $|\ell_3 = 100000|$. The level of insertion is, by default, ℓ_1 , unless stated otherwise.

For our query workloads, we consider a two-stage approach: we first identify which level our query will target according to the *levelDist* distribution; the requested value is then chosen from that level following the *valueDist* distribution. In our experiments, we use the Zipfian ($p_i \sim 1/i^\theta$) distribution for *levelDist*, while we express a bias inside each level using the uniform, 80/20, 90/10 and 99/1 distributions for *valueDist*.

Generated queries arrive at an average rate of $1 \frac{\text{query}}{\text{time_unit}}$, in almost 50k time units total simulation time. We present results for queries on a single dimension with multiple levels of hierarchy. Our default *thr* value is set to 0.3, which is a large enough value to avoid very frequent re-indexing attempts. Simulations with different values of threshold around this default show small qualitative difference. The default value of *W*, which controls how quickly the system can adapt to changes, is set to 1000 time units. For our experiments the value of *I_{max}* is set to 1000. Finally, we assume a practically infinite value of *TTL* (indices never expire).

In this section, we intend to demonstrate the performance and adaptability of our system under various conditions. Our goal is to show that we prove highly efficient under a variety of data and load distributions and can quickly adapt to sudden changes in skew without any modification to the default parameters. Specifically, we measure the percentage of queries which are answered without flooding (*precision*).

5.2 Performance Under Different Levels of Skew

In the first set of experiments we identify the behavior of our system under a variety of query loads. Specifically, we vary the number of queries directed to each level by increasing the θ parameter in the *levelDist* distribution. For each value of θ , we also choose values inside each level using four different distributions.

In Figure 5, data are skewed towards ℓ_0 . As θ increases for *levelDist*, the performance of our method improves: Re-indexing is performed sooner as more queries take place and the exact matches due to the chosen pivot level increase. By increasing the skew for *valueDist*, we observe remarkably high precision rates (close to 100%), because both the ratio of popular queries and the density of queries for certain tuples increase. Another point that plays a big role is the limited number of distinct values of ℓ_0 . Obviously this is quite small compared to the last level, thus enabling soft-indexing and faster re-indexing. For a set θ value, the method performs justifiably better as the distribution becomes more skewed: More queries exist for fewer distinct values. Finally, we notice that

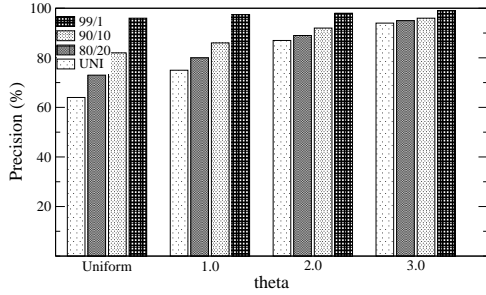


Figure 5: Precision when skew directed towards ℓ_0

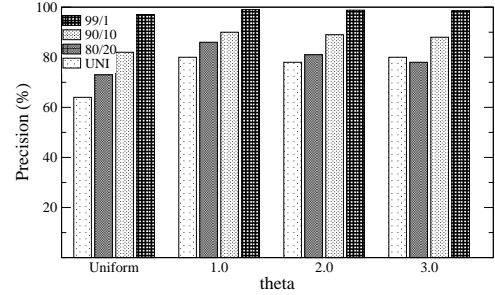


Figure 6: Precision when skew directed towards ℓ_3

the larger the θ value, the smaller the difference in precision among the different *valueDist* distributions.

Figure 6 shows results where our workload favors ℓ_3 . Again, we notice a similar trend in performance as *valueDist* becomes more biased and our method shows high precision values, albeit reduced compared to the previous case. We notice that the precision for the same *valueDist* distribution decreases as θ increases: This is due to the fact that ℓ_3 has a considerably larger number of values. By increasing the number of queries for those values, we increase (relative to the choice of *valueDist* of course) the probability of queries to non-indexed values. Nevertheless, the decrease is smaller as *valueDist* becomes more skewed.

5.3 Testing against multiple bias points

In the next experiment, we test our method against a more challenging type of workload (*MULTI*): While different levels receive an equal number of queries, nevertheless we target a *different* part of a data tree from each level. Specifically, we divide all levels in quarters and target (using different values of *valueDist*) one quarter per level so that no quarter is related with any other in the data tree. This is a very challenging workload, as it forces the method to store different data at different levels of granularity. Table 1 summarizes the results, where besides the precision we document the cost in number of re-indexed tuples as well as the number of total roll-up and drill-down operations.

Our technique proves extremely efficient in all four workloads, achieving very high precision (between 92% and 100%) at low cost: The largest number of operations occur when we uniformly query the different values in which case 75% of the tuples are re-hashed and re-inserted from re-indexing operations. As the level of skew increases, so does the number of re-indexing calls. This clearly demonstrates that our method adjusts its operation according to the need: The number of trees being re-indexed is proportional to the number of unique trees that are popular. This is a highly desirable property since for most applications we anticipate both dynamic and highly skewed loads.

5.4 Performance in dynamic environments

The adaptiveness and performance of the proposed system in a dynamic environment is examined by this set of experiments. The query distribution encloses a sudden change in skewness from level ℓ_0 towards ℓ_3 and vice versa in the middle of the simulated queries.

Figure 7 demonstrates the behavior of the system when the query load shifts from ℓ_0 towards ℓ_3 . The results show that, in all cases, our system increases its precision due to the combination of re-indexing operations and soft-state indices and the majority of questions are answered by exact match lookups. The precision reaches over 90% for $\theta = 2.0$ and over 80% for $\theta = 1.0$ before the change in skew. In the transitional stage, the flooding of the queries increases but the system rapidly manages to recover and regain its performance characteristics (after at most 5% of the queries). The

Table 1: Performance comparison for the *MULTI* workload over different values of *valueDist*

| <i>valueDist</i> | precision (%) | #roll-ups | #rolled-up trees | #drill-downs | re-inserts (%) |
|------------------|---------------|-----------|------------------|--------------|----------------|
| uniform | 92.0 | 25 | 250 | 500 | 75 |
| 80/20 | 94.3 | 25 | 250 | 171 | 42 |
| 90/10 | 95.2 | 25 | 250 | 51 | 30 |
| 99/1 | 99.5 | 1 | 10 | 6 | 1.6 |

steep decrease in precision happens at the exact time of the shift in the workload: A much larger number of distinct values belong to ℓ_3 , thus the existence of useful indices is less probable. The contribution of soft-state indices is not sufficient to handle the query load until drill-down operations take place. In this stage, the larger the value of θ , the larger the decrease in precision and the faster the recovery: As we show in Figure 8, where the query loads for *valueDist* 90/10 are considered, both exact match and indexed lookups are fewer for $\theta = 2.0$. This happens because queries are more skewed towards ℓ_3 and benefit even less from the already rolled-up trees. However, as θ increases, drill-down decisions are taken faster, favoring the increase of the exact match queries that answer the majority of the requests.

The precision of the algorithm is tested against a sudden shift from ℓ_3 to ℓ_0 for various workloads and displayed in Figure 9. During the steady stages of the simulation, similar trends to the ones of one directional skew are observed and the system presents high performance over 80% for all workloads. For more skewed *valueDist* is, the higher the precision, since the number of popular values shrinks and drill-down operations are performed faster increasing the adaptiveness of the system. After the change in the direction of skew, less queries are flooded for the $\theta = 2.0$ workloads (behavior that contrasts to the previous experiment). Figure 10 demonstrates a more comprehensive view of the system after the change in skew. Indices take over to serve lookups immediately. Due to the smaller number of distinct values in higher levels, indices perform well. However, the consecutive roll-ups destroy the existing indices and the performance of the system is influenced negatively. The system regains its performance by the rapid increase in the exact lookups.

The comparison of results among the two shifts of the workload reveals that the soft-state indices are capable to preserve the high precision of the system in case of a skew towards higher levels due to the limited number of different values. On the contrary, the adaptiveness of the system significantly depends on the re-indexing operations, when lower levels of the hierarchy are the most popular. Nevertheless, in both cases, the system needs bounded time to reorganize its indexing mechanism and achieve high performance.

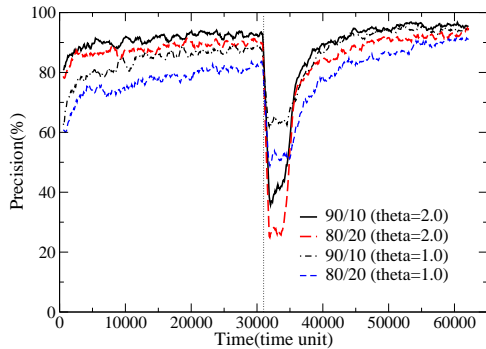


Figure 7: Precision over time for various workloads, when skewness changes from l_0 to l_3

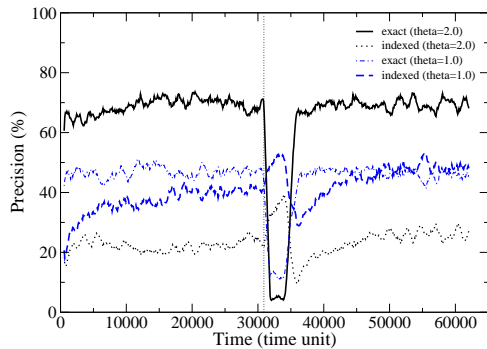


Figure 8: Precision over time of non-flooded queries for valueDist 90/10, when skew is directed from l_0 to l_3

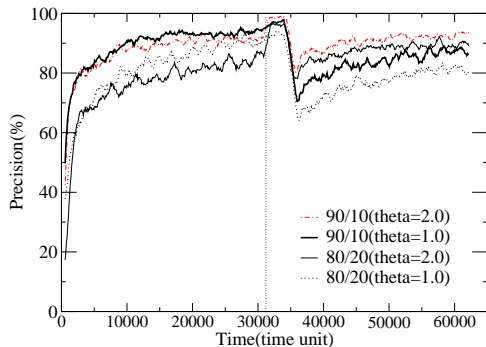


Figure 9: Precision over time for various workloads when the skewness of workload changes from l_3 to l_0

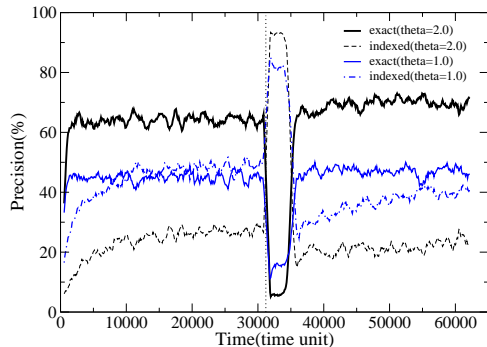


Figure 10: Precision over time of non-flooded queries for valueDist 90/10, when skew is directed from l_3 to l_0

5.5 Performance for dataset of the APB benchmark

The adaptiveness of the system is also tested using some realistic data. For this reason, we generated query sets by the APB-1 benchmark [9]. APB-1 creates a database structure with multiple dimensions and generates a set of business operations reflecting basic functionality of OLAP applications. For our experiments, we focus on the `product` dimension, a steep hierarchy of 6 levels (the bottom level contains 90% of the members). In more detail, the number of distinct values per level are $|l_0 = 50|, |l_1 = 150|, |l_2 = 800|, |l_3 = 3050|, |l_4 = 6950|$ and $|l_5 = 93050|$. Another characteristic of the specific dataset is that the number of children per node is not constant, as in the synthetically generated datasets of the previous experiments. The query load is skewed towards the lower levels of the hierarchy and 75% of queries refer to values of the l_4 and l_5 .

The results are depicted in Figure 11, where the precision over time for various initial levels of insertion is shown. It is remarkable that the system adapts to the query load and presents similar performance despite the selection of the level used as pivot level during initial insertions, thus the re-indexing operations -mainly drill-down operations towards l_4 and l_5 and soft-state indices serve to the incremental precision, which reaches values near 100%.

5.6 Updates

In order to measure the cost of incrementally updating our dataset, we randomly select the 90% of the tuples, executed each of the described query workloads in 5.2 and finally update the dataset by inserting the remaining 10% of the tuples. We note here that the workload plays an important role in the update process as it affects the indexing levels of the stored tuples and, therefore, the update cost (as tuples may have common attributes with existing

Table 2: Number of average lookups for updating indices per insertion for different values of `valueDist`

| <code>valueDist</code> | <i>Skew towards l_0</i> | | <i>Skew towards l_3</i> | |
|------------------------|--------------------------------------|----------------|--------------------------------------|----------------|
| | $\theta = 1.0$ | $\theta = 2.0$ | $\theta = 1.0$ | $\theta = 2.0$ |
| uniform | 0.0 | 0.0 | 1.99 | 1.98 |
| 80/20 | 0.01 | 0.0 | 1.8 | 1.16 |
| 90/10 | 0.14 | 0.0 | 0.39 | 0.25 |
| 99/1 | 0.0 | 0.0 | 0.01 | 0.02 |

ones). The cost in messages for storing each of the initial tuples is one lookup message so as to locate the root key and one insertion message to store the tuple. Further lookup messages are not needed, since no other indices than the ones among the root keys and corresponding pivot levels have been created yet. The selected pivot level for the initial tuples is, by default, l_1 . The conducted experiments regard the update cost in terms of additional lookups operations to inform existing indices about the appearance of the new tuples. In these set of experiments, we modified the inserted dataset. Table 2 contains the average number of lookups per insertion for updating the soft-state indices. This cost can be considered as negligible when the skew is towards high levels of the hierarchy. The maximum documented cost for skewed workloads towards l_3 and uniform `valueDist` is close to 2. The less skewed the distribution the bigger the possibility of soft-index existence in levels other than the popular one. As skew increases, this cost also diminishes.

5.7 Other experimental results

Due to space limitations, we briefly describe other conducted experiments. We experimented by varying the number of concurrent queries per time unit for the presented workloads. The experimen-

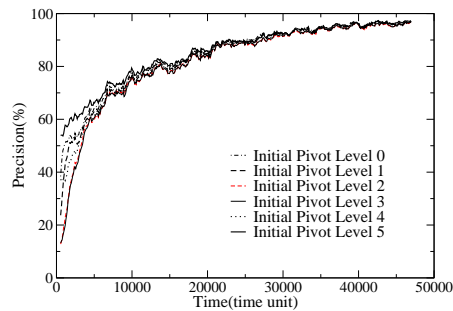


Figure 11: Precision over time for the APB workload for different initial pivot levels

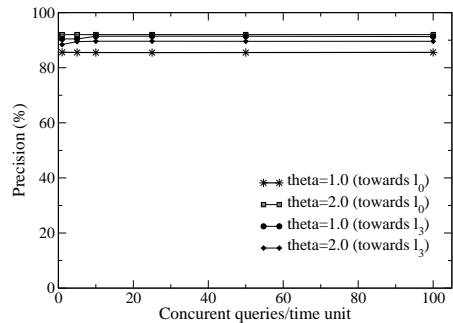


Figure 12: Precision for concurrent queries for valueDist 90/10

tal results for workloads skewed towards l_0 and l_3 are shown in Figure 12. The measured precision presents negligible variation, thus showing that the performance of the system remains consistently high while the system scales to a considerable number of concurrent users.

Moreover, experiments conducted with up to 1K nodes showed little qualitative difference. Simulation results for different values of threshold showed that the fluctuation of the precision for $0.2 \leq thr \leq 0.4$ is at most 2%. For values $thr \geq 0.5$ and uniform workloads, the fluctuation reaches 10%. For $thr \leq 0.3$, the initial number of queries to allow re-indexing should increase in order to avoid redundant operations. Experiments for various data distributions with different number of distinct values per level showed no qualitative differences. Another important observation is that by varying the default pivot level the steady-state performance of our algorithm is not affected, since the re-indexing operations and soft-state indices adapt the pivot levels appropriately.

6. CONCLUSIONS

In this work we described a highly adaptive, scalable, on-line technique in order to store hierarchical data and do efficient query processing on them. Our scheme distributes large amounts of data over a DHT overlay in a way so that the hierarchy semantics are maintained while eliminating single points of failure and minimizing data unavailability. The distinctive characteristic of our method is the adaptive indexing over the data: The stored hierarchies are indexed over variable levels according to the granularity of the incoming queries. Using limited only knowledge, peers decide on the indexing level of their stored tuples so that flooding is minimized. Moreover, there is no need for off-line updates as our system consistently updates the dataset online at low cost.

We discussed one interesting application of our method over a Grid Information System: Distributing the sources of useful data over a grid system presents significant advantages over the existing approaches. Moreover, our unique re-indexing mechanism enables

automatic aggregation of older data and more detailed views of recent ones.

Our experimental evaluation over multiple dynamic and challenging workloads confirmed our premise: Our system manages to efficiently answer the large majority of queries using very few messages. It is especially effective in skewed workloads, adapts to sudden shifts in skew and updates datasets in a fast, reliable and cost-efficient manner.

7. REFERENCES

- [1] Egee accounting portal. <http://www3.egee.cesga.es/gridsite/accounting/CESGA/>.
- [2] Ganglia Monitoring System. <http://ganglia.info/>.
- [3] GT Information Services: Monitoring and Discovery System (MDS). <http://www.globus.org/toolkit/mds/>.
- [4] Hawkeye: A Monitoring and Management Tool for Distributed Systems. <http://www.cs.wisc.edu/condor/hawkeye/>.
- [5] R-GMA: Relational Grid Monitoring Architecture. <http://www.r-gma.org/>.
- [6] The Globus Toolkit. <http://www.globus.org/>.
- [7] K. Aberer, P. Cudre-Mauroux, and M. Hauswirth. The Chatty Web: Emergent Semantics Through Gossiping. In *WWW Conference*, 2003.
- [8] K. Aberer, P. Cudre-Mauroux, M. Hauswirth, and T. V. Pelt. Gridvine: Building internet-scale semantic overlay networks. In *International Semantic Web Conference*, 2004.
- [9] *OLAP Council APB-1 OLAP Benchmark*. <http://www.olapcouncil.org/research/resrchly.htm>.
- [10] M. Ester, J. Kohlhammer, and P. Kriegel. The dc-tree: A fully dynamic index structure for data warehouses. In *ICDE*, 2000.
- [11] R. B. et.al. Apel: An implementation of grid accounting using r-gma. In *UK e-Science All Hands Conference*, 2005.
- [12] FreePastry. <http://freepastry.rice.edu/FreePastry>.
- [13] R. Huebsch, J. Hellerstein, N. L. Boon, T. Loo, S. Shenker, and I. Stoica. Querying the Internet with PIER. In *VLDB*, 2003.
- [14] V. Kantere, D. Tsoumakos, T. Sellis, and N. Roussopoulos. GrouPeer: Dynamic clustering of P2P databases. *Inf. Syst.*, 34(1):62–86, 2009.
- [15] G. Koloniari and E. Pitoura. Content-based routing of path queries in peer-to-peer systems. In *EDBT*, 2004.
- [16] L. Lakshmanan, J. Pei, and Y. Zhao. QC-trees: An Efficient Summary Structure for Semantic OLAP. In *SIGMOD*, 2003.
- [17] B. Ooi, Y. Shu, K. Tan, and A. Zhou. PeerDB: A P2P-based System for Distributed Data Sharing. In *ICDE*, 2003.
- [18] Y. Sismanis, A. Deligiannakis, Y. Kotidis, and N. Roussopoulos. Hierarchical dwarfs for the rollup cube. In *DOLAP*, 2003.
- [19] C. Tang, Z. Xu, and S. Dwarkadas. Peer-to-peer information retrieval using self-organizing semantic overlay networks. In *SIGCOMM*, 2003.
- [20] I. Tatarinov and A. Halevy. Efficient Query Reformulation in Peer-Data Management Systems. In *SIGMOD*, 2004.
- [21] W. Wang, H. Lu, J. Feng, and J. X. Yu. Condensed Cube: An Effective Approach to Reducing Data Cube Size. In *ICDE*, 2002.
- [22] X. Zhang, J. Freschl, and J. Schopf. Scalability analysis of three monitoring and information systems: MDS2, R-GMA, and Hawkeye. *J. Parallel Distrib. Comput.*, 67(8):883–902, 2007.