

The design and architecture of the τ -Synopsis system

Yossi Matias, Leon Portman, Natasha Drukh

School of Computer Science

Tel Aviv University

matias@cs.tau.ac.il, leon.portman@nice.com, kreimern@cs.tau.ac.il

Abstract. Data synopses are concise representations of data sets, that enable effective processing of approximate queries to the data sets. The τ -Synopsis is a system designed to provide a run-time environment for remote execution of multiple synopses for both relational as well as XML databases. The system can serve as an effective research platform for experimental evaluation and comparison of different synopses, as well as a platform for studying the effective management of multiple synopses in a federated or centralized environment.

1 Introduction

In large data recording and warehousing environments, it is often advantageous to provide fast, approximate answers to queries, whenever possible. The goal is to provide a quick response in orders of magnitude faster than the time to compute an exact answer, by avoiding or minimizing the number of accesses to the base data.

Approximate query processing is supported by synopses that are compact representations of the original data, such as histograms, samples, wavelet-synopses or other methods [1]. In the AQUA system [2], synopses are precomputed and stored in a DBMS. The system supports approximate answers by rewriting queries originally directed to the base tables to run on these synopses, and it enables keeping synopses up-to-date as the database changes. The question of how to reconcile various synopses for large data sources with many tables was studied in [3].

The τ -Synopsis system was designed to provide a run-time environment for execution of multiple synopses. The system can serve as an effective research platform for experimental evaluation and comparison of different synopses, as well as a platform for studying the effective management of multiple synopses. The synopses can be placed at a centralized environment, or they can function as web services in a federated architecture.

A software demo of the system as a federated environment with remote execution of synopses was presented in [5]. A software demo of the system with emphasis on the synopses management in a centralized environment was presented in [4]. The system currently includes several dozens of synopses for both Relational as well as XML databases.

This paper presents a high-level overview of the architecture and functionality of the system. For more details, please refer to the full paper [6].

2 The τ -Synopsis Functionality

The main operational processes supported by the τ -Synopsis system are: constructing and updating multiple pluggable synopses, interception and analysis of query workload, interception and analysis of data updates, approximate query processing, synopses management, and benchmarking.

The user interface provides an administrator user with a capability to manage data sources, synopses specifications, updates and pre-defined workloads. Figure 1 depicts the main administration UI.

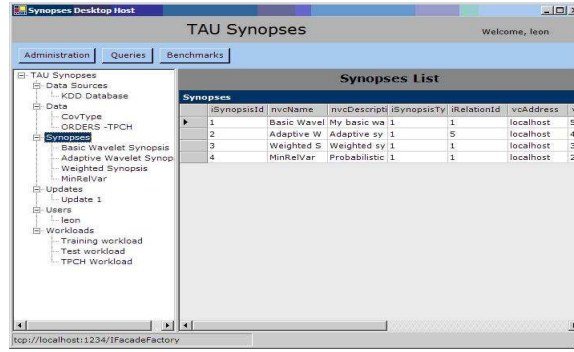


Fig. 1. The Administration UI

End users can test and compare different synopses that are registered in the system. In the *Query execution* mode a user can evaluate a single synopsis at a time.

In the *Query Mode*, the user selects the synopsis to be evaluated. Relational queries can be of the following structure:

```
SELECT Sum(Data) FROM Relation WHERE filter > l AND filter < h .
```

For XML synopses, the queries are XPath expressions. The system validates the user input expression for the XPath syntax and the tag labels are validated against existing labels of the underlying XML document.

The Query mode also allows the evaluation of multiple queries at a time by specifying a workload to be evaluated. The approximate results obtained using the registered synopses are depicted together with the exact results computed by the system.

The *Benchmark Mode* enables multiple synopses evaluation over pre-defined workloads and their comparison using visual display; see Figure 2. The user selects the synopses to be evaluated and the workload to be used for the evaluation. For the performance measurements, the minimum, maximum and step size

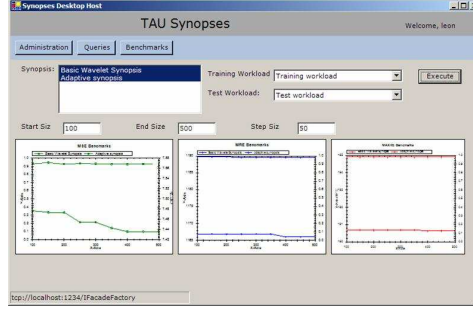


Fig. 2. Benchmark Mode

for the synopses construction are set by the user. The system invokes the construction of the different synopses, and these synopses are then evaluated over the selected workload. The system can compute the accuracy of the different synopses using several error metrics.

3 Architecture

In order to provide an effective operational and research platform the τ -Synopses system commits to the following design goals:

- *Pluggable integration*
- *Remote execution*
- *Distributed client-server environment*
- *Flexibility and scalability*
- *Low bandwidth requirement*

The core of the τ -Synopses system architecture features the following components: Query Execution Engine, Synopses Manager, Updates Logger, and Workload Manager. These modules interact with a relational or XML databases which hold the data sets, and with registered synopses that act as web services. These synopses are connected either locally or remotely through a SOAP-enabled platforms.

The Synopses Manager is used for registration and maintenance of the synopses. A new synopsis is added to the system by registering its parameters (including list of supported queries and data sets) in the Synopses Manager Catalog.

The Query Execution Engine supports an interface for receiving query request from end-users and invoking the appropriate synopsis (or synopses), as determined by the Synopses Manager in order to process such query.

The Updates Logger feeds all data updates to the registered synopses by intercepting data updates information in the data sources.

The Workload Manager captures, maintains and analyzes workload information for building, maintaining and testing synopses.

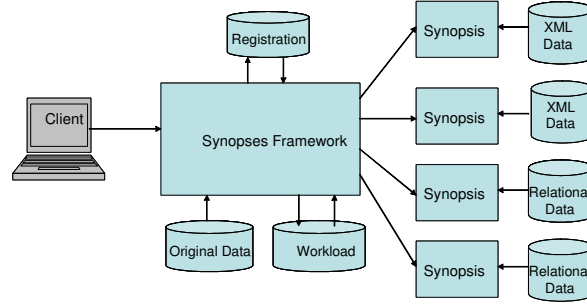


Fig. 3. General Architecture

The system provides a light-weight host process, inside which the custom synopses will be running. The host is responsible for all communication with the system and is transparent to the synopsis. This design enables unconstrained deployment. A remote synopsis can be integrated into the system by deploying or adapting such host into the remote system, and connecting the synopsis module locally into the host. Figure 3 illustrates an overall view of the system in a distributed environment, consisting of multiple remote synopses, each representing its local data source.

The system modules were implemented in the .NET framework, with remote modules communicating through the .NET Remoting. Any relational DB can be used as a database provider.

Acknowledgement. We thank Yariv Matia and Daniel Urieli for their contributions to the system development. We also thank the students in various classes at Tel Aviv university who have contributed synopses implementations to the system and for their feedback.

References

1. P. B. Gibbons and Y. Matias. Synopsis data structures for massive data sets. *External Memory Algorithms, DIMACS Series in Discrete Mathematics and Theoretical Computer Science, American Mathematical Society*, 50, A, 1999. Also in SODA'99.
2. P. B. Gibbons, V. Poosala, S. Acharya, Y. Bartal, Y. Matias, S. Muthukrishnan, S. Ramaswamy, and T. Suel. AQUA: System and techniques for approximate query answering. Technical report, Bell Laboratories, Murray Hill, 1998.
3. A. C. König and G. Weikum. A framework for the physical design problem for data synopses. In *EDBT 2002 - Advances in Database Technology*, March 2002.
4. Y. Matia, Y. Matias, and L. Portman. Synopses reconciliation via calibration in the τ -synopses system. In *Software Demo, EDBT'06*.
5. Y. Matias and L. Portman. τ -Synopses: a system for run-time management of remote synopses. In *Software Demo, ICDE'04, EDBT'04*.
6. Y. Matias, L. Portman, and N. Drukh. The design and architecture of the τ -synopses system. Technical Report, Tel Aviv University, 2005.