



1 Stream Cube: An Architecture 2 for Multi-Dimensional Analysis of Data Streams

3 JIAWEI HAN

hanj@cs.uiuc.edu

4 *University of Illinois*

5 YIXIN CHEN

chen@cse.wustl.edu

6 *Washington University, St. Louis*

7 GUOZHU DONG

gdong@cs.wright.edu

8 *Wright State University*

9 JIAN PEI

jpei@cs.sfu.ca

10 *Simon Fraser University, B. C., Canada*

11 BENJAMIN W. WAH

b-wah@uiuc.edu

12 *University of Illinois*

13 JIANYONG WANG

jianyong@tsinghua.edu.cn

14 *Tsinghua University, Beijing, China*

15 Y. DORA CAI

ycai@ncsa.uiuc.edu

16 *University of Illinois*

17 **Recommended by:** Ahmed Elmagarmid

18 **Published online:** xx xx

19 **Abstract.** Real-time surveillance systems, telecommunication systems, and other dynamic environments often
20 generate tremendous (potentially infinite) volume of stream data: the volume is too huge to be scanned multiple
21 times. Much of such data resides at rather low level of abstraction, whereas most analysts are interested in relatively
22 high-level dynamic changes (such as trends and outliers). To discover such high-level characteristics, one may need
23 to perform on-line multi-level, multi-dimensional analytical processing of stream data. In this paper, we propose
24 an architecture, called *stream_cube*, to facilitate on-line, multi-dimensional, multi-level analysis of stream data.

25 For fast online multi-dimensional analysis of stream data, three important techniques are proposed for efficient
26 and effective computation of stream cubes. First, a *tilted time frame* model is proposed as a multi-resolution model
27 to register time-related data: the more recent data are registered at finer resolution, whereas the more distant data
28 are registered at coarser resolution. This design reduces the overall storage of time-related data and adapts nicely
29 to the data analysis tasks commonly encountered in practice. Second, instead of materializing cuboids at all levels,
30 we propose to maintain a small number of *critical* layers. Flexible analysis can be efficiently performed based on
31 the concept of *observation layer* and *minimal interesting layer*. Third, an efficient stream data cubing algorithm
32 is developed which computes only the layers (cuboids) along a *popular path* and leaves the other cuboids for
33 query-driven, on-line computation. Based on this design methodology, stream data cube can be constructed and
34 maintained incrementally with a reasonable amount of memory, computation cost, and query response time. This
35 is verified by our substantial performance study.

Stream data cube architecture facilitates online analytical processing of stream data. It also forms a preliminary data structure for online stream data mining. The impact of the design and implementation of stream data cube in the context of stream data mining is also discussed in the paper.

1. Introduction

With years of research and development of data warehousing and OLAP technology [9, 15], a large number of data warehouses and data cubes have been successfully constructed and deployed in applications, and data cube has become an essential component in most data warehouse systems and in some extended relational database systems and has been playing an increasingly important role in data analysis and intelligent decision support.

The data warehouse and OLAP technology is based on the integration and consolidation of data in multi-dimensional space to facilitate powerful and fast on-line data analysis. Data are aggregated either completely or partially in multiple dimensions and multiple levels, and are stored in the form of either relations or multi-dimensional arrays [1, 29]. The dimensions in a data cube are of categorical data, such as products, region, time, etc., and the measures are numerical data, representing various kinds of aggregates, such as *sum*, *average*, *variance* of sales or profits, etc.

The success of OLAP technology naturally leads to its possible extension from the analysis of static, pre-integrated, historical data to that of current, dynamically changing data, including time-series data, scientific and engineering data, and data produced in other dynamic environments, such as power supply, network traffic, stock exchange, telecommunication data flow, Web click streams, weather or environment monitoring, etc.

A fundamental difference in the analysis of stream data from that of relational and warehouse data is that the stream data is generated in huge volume, flowing in-and-out dynamically, and changing rapidly. Due to limited memory or disk space and processing power available in today's computers, most data streams may only be examined in a single pass. These characteristics of stream data have been emphasized and investigated by many researchers, such as [6, 7, 12, 14, 16], and efficient stream data querying, clustering and classification algorithms have been proposed recently (such as [12, 14, 16, 17, 20]). However, there is another important characteristic of stream data that has not drawn enough attention: *Most of stream data resides at rather low level of abstraction, whereas an analyst is often more interested in higher and multiple levels of abstraction.* Similar to OLAP analysis of static data, multi-level, multi-dimensional on-line analysis should be performed on stream data as well.

The requirement for multi-level, multi-dimensional on-line analysis of stream data, though desirable, raises a challenging research issue: *"Is it feasible to perform OLAP analysis on huge volumes of stream data since a data cube is usually much bigger than the original data set, and its construction may take multiple database scans?"*

In this paper, we examine this issue and present an interesting architecture for on-line analytical analysis of stream data. Stream data is generated continuously in a dynamic environment, with huge volume, infinite flow, and fast changing behavior. As collected, such data is almost always at rather low level, consisting of various kinds of detailed temporal and other features. To find interesting or unusual patterns, it is essential to perform analysis on some useful measures, such as sum, average, or even more sophisticated measures, such

79 as regression, at certain meaningful abstraction level, discover critical changes of data, and
80 drill down to some more detailed levels for in-depth analysis, when needed.

81 To illustrate our motivation, let's examine the following examples.

82 *Example 1.* A power supply station can watch infinite streams of power usage data, with
83 the lowest granularity as individual user, location, and second. Given a large number of
84 users, it is only realistic to analyze the fluctuation of power usage at certain high levels,
85 such as by city or street district and by quarter (of an hour), making timely power supply
86 adjustments and handling unusual situations.

87 Conceptually, for multi-dimensional analysis, one can view such stream data as a *virtual*
88 data cube, consisting of one or a few measures and a set of dimensions, including one *time*
89 *dimension*, and a few other dimensions, such as location, user-category, etc. However, in
90 practice, it is impossible to materialize such a data cube, since the materialization requires a
91 huge amount of data to be computed and stored. Some efficient methods must be developed
92 for systematic analysis of such data.

93 *Example 2.* Suppose that a Web server, such as Yahoo.com, receives a huge volume of
94 Web click streams requesting various kinds of services and information. Usually, such
95 stream data resides at rather low level, consisting of time (down to subseconds), Web page
96 address (down to concrete URL), user ip address (down to detailed machine IP address),
97 etc. However, an analyst may often be interested in changes, trends, and unusual patterns,
98 happening in the data streams, at certain high levels of abstraction. For example, it is
99 interesting to find that *the Web clicking traffic in North America on sports in the last*
100 *15 minutes is 40% higher than the last 24 hours' average.*

101 From the point of view of a Web analysis provider, given a large volume of fast changing
102 Web click streams, and with limited resource and computational power, it is only realistic
103 to analyze the changes of Web usage at certain high levels, discover unusual situations,
104 and drill down to some more detailed levels for in-depth analysis, when needed, in order to
105 make timely responses.

106 Interestingly, both the analyst and analysis provider share a similar view on such stream
107 data analysis: instead of bogging down to every detail of data stream, a demanding request is
108 to provide on-line analysis of changes, trends and other patterns at high levels of abstraction,
109 with low cost and fast response time.

110 In this study, we take Example 2 as a typical scenario and study how to perform efficient
111 and effective multi-dimensional analysis of stream data, with the following contributions.

- 112 1. For on-line stream data analysis, both space and time are critical. In order to avoid im-
113 posing unrealistic demand on space and time, instead of computing a fully materialized
114 cube, we suggest to compute a partially materialized data cube, with a *tilted time frame*
115 as its time dimension model. In the *tilted time frame*, time is registered at different levels
116 of granularity. The most recent time is registered at the finest granularity; the more
117 distant time is registered at coarser granularity; the level of coarseness depends on the
118 application requirements and on how old the time point is. This model is sufficient for
119 most analysis tasks, and at the same time it also ensures that the total amount of data to
120 retain in memory or to be stored on disk is small.

2. Due to limited memory space in stream data analysis, it is often too costly to store a precomputed cube, even with the *tilted time frame*, which substantially compresses the storage space. We propose to compute and store only two *critical layers* (which are essentially cuboids) in the cube: (1) an *observation layer*, called *o-layer*, which is the layer that an analyst would like to check and make decisions for either signaling the exceptions or drilling on the exception cells down to lower layers to find their corresponding lower level exceptions; and (2) the *minimal interesting layer*, called *m-layer*, which is the minimal layer that an analyst would like to examine, since it is often neither cost-effective nor practically interesting to examine the minute detail of stream data. For example, in Example 1, we assume that the *o-layer* is *user-region*, *theme*, and *quarter*, while the *m-layer* is *user*, *sub-theme*, and *minute*.
3. Storing a cube at only two critical layers leaves a lot of room at what to compute and how to compute for the cuboids between the two layers. We propose one method, called popular-path cubing, which rolls up the cuboids from the *m-layer* to the *o-layer*, by following one popular drilling path, materializes only the layers along the path, and leave other layers to be computed only when needed. Our performance study shows that this method achieves a reasonable trade-off between space, computation time, and flexibility, and has both quick aggregation time and exception detection time.

The rest of the paper is organized as follows. In Section 2, we define the basic concepts and introduce the research problem. In Section 3, we present an architectural design for online analysis of stream data by defining the problem and introducing the concepts of *tilted time frame* and *critical layers*. In Section 4, we present the *popular-path* cubing method, an efficient algorithm for stream data cube computation that supports on-line analytical processing of stream data. Our experiments and performance study of the proposed methods are presented in Section 5. The related work and possible extensions of the model are discussed in Section 6, and our study is concluded in Section 7.

2. Problem definition

In this section, we introduce the basic concepts related to data cubes, multi-dimensional analysis of stream data, and stream data cubes, and define the problem of research.

The concept of data cube [15] was introduced to facilitate multi-dimensional, multi-level analysis of large data sets.

Let \mathcal{D} be a relational table, called the base table, of a given cube. The set of all *attributes* \mathcal{A} in \mathcal{D} are partitioned into two subsets, the *dimensional attributes* DIM and the *measure attributes* M (so $DIM \cup M = \mathcal{A}$ and $DIM \cap M = \emptyset$). The measure attributes functionally depend on the dimensional attributes in \mathcal{DB} and are defined in the context of data cube using some typical aggregate functions, such as COUNT, SUM, AVG, or some more sophisticated computational functions, such as standard deviation, regression, etc.

A tuple with schema \mathcal{A} in a multi-dimensional space (i.e., in the context of data cube) is called a **cell**. Given three distinct cells c_1 , c_2 and c_3 , c_1 is an **ancestor** of c_2 , and c_2 a **descendant** of c_1 iff on every dimensional attribute, either c_1 and c_2 share the same value, or c_1 's value is a generalized value of c_2 's in the dimension's concept hierarchy. c_2 is a

162 **sibling** of c_3 iff c_2 and c_3 have identical values in all dimensions except one dimension A
 163 where $c_2[A]$ and $c_3[A]$ have the same parent in the dimension's domain hierarchy. A cell
 164 which has k non-* values is called a k -d **cell**. (We use "*" to indicate "all", i.e., the highest
 165 level on any dimension.)

166 A tuple $c \in \mathcal{D}$ is called a **base cell**. A base cell does not have any descendant. A cell c
 167 is an **aggregated cell** iff it is an ancestor of some base cell. For each aggregated cell c , its
 168 values on the measure attributes are derived from the complete set of descendant base cells
 169 of c . An aggregated cell c is an **iceberg cell** iff its measure value satisfies a specified iceberg
 170 condition, such as $\text{measure} \geq \text{val}_1$. The data cube that consists of all and only the iceberg
 171 cells satisfying a specified iceberg condition I is called the **iceberg cube** of a database \mathcal{D}
 172 under condition I .

173 Notice that in stream data analysis, besides the popularly used SQL aggregate-based
 174 measures, such as COUNT, SUM, MAX, MIN, and AVG, *regression* is a useful measure.
 175 A stream data cell compression technique LCR (*linearly compressed representation*) is
 176 developed in [10] to support efficient on-line regression analysis of stream data in data
 177 cubes. The study [10] shows that for linear and multiple linear regression analysis, only a
 178 small number of *regression measures* rather than the complete stream of data need to be
 179 registered. This holds for regression on both the time dimension and the other dimensions.
 180 Since it takes a much smaller amount of space and time to handle regression measures in
 181 a multi-dimensional space than handling the stream data itself, it is preferable to construct
 182 regression (-measured) cubes by computing such regression measures.

183 A *data stream* is considered as a huge volume, infinite flow of data records, such as Web
 184 click streams, telephone call logs, and on-line transactions. The data is collected at the
 185 most detailed level in a multi-dimensional space, which may represent time, location, user,
 186 theme, and other semantic information. Due to the huge amount of data and the transient
 187 behavior of data streams, most of the computations will scan a data stream only once.
 188 Moreover, the direct computation of measures at the most detailed level may generate a
 189 huge number of results but may not be able to disclose the general characteristics and
 190 trends of data streams. Thus data stream analysis will require to consider aggregations and
 191 analysis at multi-dimensional and multi-level space.

192 Our task is to support *efficient, high-level, on-line, multi-dimensional analysis of such*
 193 *data streams in order to find unusual (exceptional) changes of trends, according to users'*
 194 *interest, based on multi-dimensional numerical measures*. This may involve construction
 195 of a data cube, if feasible, to facilitate on-line, flexible analysis.

196 3. Architecture for on-line analysis of data streams

197 To facilitate on-line, multi-dimensional analysis of data streams, we propose a stream_cube
 198 architecture with the following features: (1) *tilted time frame*, (2) two *critical layers*: a
 199 *minimal interesting layer* and an *observation layer*, and (3) *partial computation of data*
 200 *cubes by popular-path cubing*. The stream data cubes so constructed are much smaller than
 201 those constructed from the raw stream data but will still be effective for multi-dimensional
 202 stream data analysis tasks.

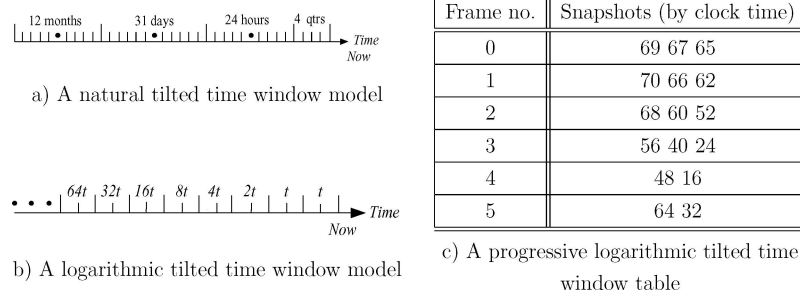


Figure 1. Three models for tilted time windows.

203 3.1. Tilted time frame

204 In stream data analysis, people are usually interested in recent changes at a fine scale, but
 205 long term changes at a coarse scale. Naturally, one can register time at different levels of
 206 granularity. The most recent time is registered at the finest granularity; the more distant time
 207 is registered at coarser granularity; and the level of coarseness depends on the application
 208 requirements and on how old the time point is (from the current time).

209 There are many possible ways to design a titled time frame. We adopt three kinds of
 210 models: (1) *natural tilted time window model* (figure 1(a)), (2) *logarithmic scale tilted*
 211 *time window model* (figure 1(b)), and (3) *progressive logarithmic tilted time window model*
 212 (figure 1(c)).

213 A *natural tilted time window model* is shown in figure 1(a), where the time frame is
 214 structured in multiple granularity based on natural time scale: the most recent 4 quarters
 215 (15 minutes), then the last 24 hours, 31 days, and 12 months (the concrete scale will be
 216 determined by applications). Based on this model, one can compute frequent itemsets in
 217 the last hour with the precision of quarter of an hour, the last day with the precision of hour,
 218 and so on, until the whole year, with the precision of month.¹ This model registers only $4 +$
 219 $24 + 31 + 12 = 71$ units of time for a year instead of $366 \times 24 \times 4 = 35,136$ units, a saving
 220 of about 495 times, with an acceptable trade-off of the grain of granularity at a distant time.

221 The second choice is *logarithmic tilted time model* as shown in figure 1(b), where the time
 222 frame is structured in multiple granularity according to a logarithmic scale. Suppose the
 223 current window holds the transactions in the current quarter. Then the remaining slots are
 224 for the last quarter, the next two quarters, 4 quarters, 8 quarters, 16 quarters, etc., growing at
 225 an exponential rate. According to this model, with one year of data and the finest precision
 226 at quarter, we will need $\lceil \log_2(365 \times 24 \times 4) + 1 \rceil = 17$ units of time instead of $366 \times$
 227 $24 \times 4 = 35,136$ units. That is, we will just need 17 time frames to store the compressed
 228 information.

229 The third choice is a *progressive logarithmic tilted time frame*, where snapshots are stored
 230 at differing levels of granularity depending upon the recency. Snapshots are classified into
 231 different *frame number* which can vary from 1 to *max.frame*, where $\log_2(T) - \text{max-capacity}$
 232 $\leq \text{max.frame} \leq \log_2(T)$, *max-capacity* is the maximal number of snapshots held in each
 233 frame, and T is the clock time elapsed since the beginning of the stream.

Each snapshot is represented by its timestamp. The rules for insertion of a snapshot t (at time t) into the snapshot frame table are defined as follows: (1) if $(t \bmod 2^i) = 0$ but $(t \bmod 2^{i+1}) \neq 0$, t is inserted into *frame_number* i if $i \leq \text{max_frame}$; otherwise (i.e., $i > \text{max_frame}$), t is inserted into *max_frame*; and (2) each slot has a *max_capacity* (which is 3 in our example of figure 1(c)). At the insertion of t into *frame_number* i , if the slot already reaches its *max_capacity*, the oldest snapshot in this frame is removed and the new snapshot is inserted. For example, at time 70, since $(70 \bmod 2^1) = 0$ but $(70 \bmod 2^2) \neq 0$, 70 is inserted into frame-number 1 which knocks out the oldest snapshot 58 if the slot capacity is 3. Also, at time 64, since $(64 \bmod 2^6) = 0$ but $\text{max_frame} = 5$, so 64 has to be inserted into frame 5. Following this rule, when slot capacity is 3, the following snapshots are stored in the tilted time window table: 16, 24, 32, 40, 48, 52, 56, 60, 62, 64, 65, 66, 67, 68, 69, 70, as shown in figure 1(c). From the table, one can see that the closer to the current time, the denser are the snapshots stored.

In the logarithmic and progressive logarithmic models discussed above, we have assumed that the base is 2. Similar rules can be applied to any base α , where α is an integer and $\alpha > 1$. The tilted time models shown above are sufficient for usual time-related queries, and at the same time it ensures that the total amount of data to retain in memory and/or to be computed is small.

Both the natural tilted window model and the progressive logarithmic tilted time window model provide a natural and systematic way for incremental insertion of data in new windows and gradually fading out the old ones. To simplify our discussion, we will only use the natural titled time window model in the following discussions. The methods derived from this time window can be extended either directly or with minor modifications to other time windows.

In our data cube design, we assume that each cell in the base cuboid and in an aggregate cuboid contains a tilted time frame, for storing and propagating measures in the computation. This tilted time window model is sufficient to handle usual time-related queries and mining, and at the same time it ensures that the total amount of data to retain in memory and/or to be computed is small.

3.2. Critical layers

Even with the *tilted time frame* model, it could still be too costly to dynamically compute and store a full cube since such a cube may have quite a few dimensions, each containing multiple levels with many distinct values. Since stream data analysis has only limited memory space but requires fast response time, a realistic arrangement is to compute and store only some mission-critical cuboids in the cube.

In our design, two critical cuboids are identified due to their conceptual and computational importance in stream data analysis. We call these cuboids layers and suggest to compute and store them dynamically. The first layer, called *m-layer*, is the minimally interesting layer that an analyst would like to study. It is necessary to have such a layer since it is often neither cost-effective nor practically interesting to examine the minute detail of stream data. The second layer, called *o-layer*, is the observation layer at which an analyst (or an automated system) would like to check and make decisions of either signaling the

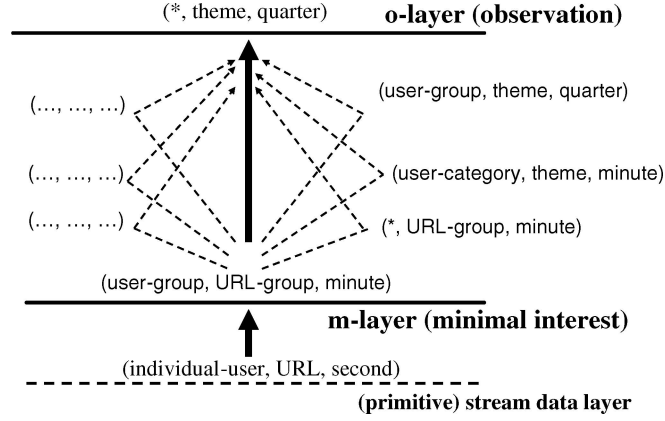


Figure 2. Two critical layers in the stream cube.

276 exceptions, or drilling on the exception cells down to lower layers to find their lower-level
 277 exceptional descendants.

278 *Example 3.* Assume that “(individual-user, URL, second)” forms the primitive layer of the
 279 input stream data in Example 1. With the *tilted time frame* as shown in figure 1, the two
 280 critical layers for power supply analysis are: (1) the *m*-layer: (user-group, URL-group,
 281 minute), and (2) the *o*-layer: (*, theme, quarter), as shown in figure 2.

282 Based on this design, the cuboids lower than the *m*-layer will not need to be computed
 283 since they are out of the minimal interest of users. Thus the minimal interesting cells that
 284 our base cuboid needs to compute and store will be the aggregate cells computed with
 285 grouping by user-group, URL-group, and minute. This can be done by aggregations (1) on
 286 two dimensions, user and URL, by rolling up from individual-user to user-group and from
 287 URL to URL-group, respectively, and (2) on time dimension by rolling up from second to
 288 minute.

289 Similarly, the cuboids at the *o*-layer should be computed dynamically according to the
 290 tilted time frame model as well. This is the layer that an analyst takes as an observation
 291 deck, watching the changes of the current stream data by examining the slope of changes
 292 at this layer to make decisions. The layer can be obtained by rolling up the cube (1) along
 293 two dimensions to * (which means all user-category) and theme, respectively, and (2) along
 294 time dimension to quarter. If something unusual is observed, the analyst can drill down to
 295 examine the details and the exceptional cells at low levels.

296 3.3. Partial materialization of stream cube

297 Materializing a cube at only two critical layers leaves much room for how to compute the
 298 cuboids in between. These cuboids can be precomputed fully, partially, not at all (i.e., leave
 299 everything computed on-the-fly), or precomputing exception cells only. Let us first examine
 300 the feasibility of each possible choice in the environment of stream data. Since there may

be a large number of cuboids between these two layers and each may contain many cells, it is often too costly in both space and time to fully materialize these cuboids, especially for stream data. Moreover, for the choice of computing *exception cells* only, the problem becomes how to set up an exception threshold. A too low threshold may lead to computing almost the whole cube, whereas a too high threshold may leave a lot of cells uncomputed and thus not being able to answer many interesting queries efficiently. On the other hand, materializing nothing forces all the aggregate cells to be computed on-the-fly, which may slow down the response time substantially. Thus, it seems that the only viable choice is to perform partial materialization of a stream cube.

According to the above discussion, we propose the following framework in our computation.

Framework 3.1 (Partial materialization of stream data). The task of computing a stream data cube is to (1) compute two critical layers (cuboids): (i) *m-layer* (the minimal interest layer), and (ii) *o-layer* (the observation layer), and (2) materialize only a reasonable fraction of the cuboids between the two layers which can allow efficient on-line computation of other cuboids.

Partial materialization of data cubes has been studied in previous works [9, 19]. With the concern of both space and on-line computation time, the partial computation of stream data cube poses more challenging issues than its static counterpart: partial computation of nonstream data cubes, since we have to ensure not only the limited size of the precomputed cube and limited precomputation time, but also efficient online incremental updating upon the arrival of new stream data, as well as fast online drilling to find interesting aggregates and patterns. Obviously, such partial computation should lead to the computation of a rather small number of cuboids, fast updating, and fast online drilling. We will examine how to design such a stream data cube in the next section.

4. Stream data cube computation

From the above analysis, one can see that in order to design an efficient and scalable stream data cube, it is essential to lay out clear design requirements so that we can ensure that the cube can be computed and maintained efficiently in the stream data environment and can provide fast online multidimensional stream data analysis. We have the following design requirements.

1. A stream data cube should be relatively stable in size with respect to infinite data streams. Since a stream data cube takes a set of potentially infinite data streams as inputs, if the size of the base-cuboid grows indefinitely with the size of data streams, the size of stream data cube will grow indefinitely. It is impossible to realize such a stream data cube. Fortunately, with tilted time frame, the distant time is compressed substantially and the very distant data beyond the specified time frame are faded out (i.e., removed) according to the design. Thus the bounded time frames transform infinite data streams into finite, compressed representation, and if the data in the other dimensions of the base cuboid are relatively stable with time, the entire base-cuboid (with the time dimensions included) should be relatively stable in size.

2. A stream data cube should be incrementally updateable with respect to infinite data streams. Since a stream data cube takes potentially infinite data streams as inputs, it is impossible to construct the cube from scratch and the cube must be incrementally updatable. Any cube design that is not incrementally updatable cannot be used as the architecture of a stream cube.
 3. The time taken for incremental computation of a stream data cube should be proportional to the size of the incremental portion of the base cuboid of the cube. To incrementally update a stream data cube, one must start from the incremental portion of the base cuboid and use an efficient algorithm to compute it. The time to compute such an incremental portion of the cube should be proportional (desirably, linear) to the size of the incremental portion of the base cuboid of the cube.
 4. The stream data cube should facilitate the fast online drilling along any single dimension or along the combination of a small number of dimensions. Although it is impossible to materialize all the cells of a stream cube, it is expected that the drilling along a single dimension or along the combination of a small number of dimensions be fast. Materialization of some portion of the cube will facilitate such fast online presentation.
- Based on the above design requirements, we examine the methods for the efficient computation of stream cubes.

4.1. Design of stream cube architecture: A popular path architecture

According to our discussion in Section 3, there are three essential components in a stream data cube: (1) *tilted time frame*, (2) *two critical layers: a minimal interesting layer and an observation layer*, and (3) *partial computation of data cubes*.

In data cube computation, iceberg cube [8] which stores only the aggregate cells that satisfy an iceberg condition has been used popularly as a data cube architecture since it may substantially reduce the size of a data cube when data is sparse. In stream data analysis, people may often be interested in only the substantially important or exceptional cube cells, and such important or exceptional conditions can be formulated as typical *iceberg conditions*. Thus it seems that iceberg cube could be an interesting model for stream cube architecture. Unfortunately, iceberg cube cannot accommodate the incremental update with the constant arrival of new data and thus cannot be used as the architecture of stream data cube. We have the following observation.

Framework 4.1 (No iceberg cubing for stream data). The iceberg cube model does not fit the stream data cube architecture. Nor does the exceptional cube model.

Rationale. With the incremental and gradual arrival of new stream data, as well as the incremental fading of the obsolete data from the time scope of a data cube, it is required that incremental update be performed on such a stream data cube. It is unrealistic to constantly recompute the data cube from scratch upon incremental updates due to the tremendous cost of recomputing the cube on the fly. Unfortunately, such an incremental model does not fit the iceberg cube computation model due to the following observation: Let a cell " $\langle d_i, \dots, d_k \rangle: m_{ik}$ " represent a $k - i + 1$ dimension cell with d_i, \dots, d_k as its corresponding dimension values and m_{ik} as its measure value. If $SAT(m_{ik}, iceberg_cond)$ is

383 false, i.e., m_{ik} does not satisfy the iceberg condition, the cell is dropped from the iceberg
 384 cube. However, at a later time slot t' , the corresponding cube cell may get a new measure
 385 m'_{ik} related to t' . Since m_{ik} has been dropped at a previous instance of time due to its
 386 inability to satisfy the iceberg condition, the new measure for this cell cannot be calculated
 387 correctly without such information. Thus one cannot use the iceberg architecture to model a
 388 stream data cube unless recomputing the measure from the based cuboid upon each update.
 389 Similar reasoning can be applied to the case of exceptional cell cubes since the exceptional
 390 condition can be viewed as a special iceberg condition.

391 Since iceberg cube cannot be used as a stream cube model, but materializing the full
 392 cube is too costly both in computation time and storage space, we propose to compute only
 393 a *popular path* of the cube as our partial computation of stream data cube, as described
 394 below.

395 Based on the notions of the minimal interesting layer (the m -layer) and the tilted time
 396 frame, stream data can be directly aggregated to this layer according to the tilted time
 397 scale. Then the data can be further aggregated following one popular drilling path to
 398 reach the observation layer. That is, the *popular path* approach computes and maintains a
 399 single popular aggregation path from m -layer to o -layer so that queries directly on those
 400 (layers) along the popular path can be answered without further computation, whereas
 401 those deviating from the path can be answered with minimal online computation from
 402 those reachable from the computed layers. Such cost reduction makes possible the OLAP-
 403 styled exploration of cubes in stream data analysis.

404 To facilitate efficient computation and storage of the popular path of the stream cube,
 405 a compact data structure needs to be introduced so that the space taken in the compu-
 406 tation of aggregations is minimized. A data structure, called H-tree, a hyper-linked tree
 407 structure introduced in [18], is revised and adopted here to ensure that a compact structure
 408 is maintained in memory for efficient computation of multi-dimensional and multi-level
 409 aggregations.

410 We present these ideas using an example.

411 *Example 4.* Suppose the stream data to be analyzed contains 3 dimensions, A, B and C , each
 412 with 3 levels of abstraction (excluding the highest level of abstraction “*”), as (A_1, A_2, A_3) ,
 413 (B_1, B_2, B_3) , (C_1, C_2, C_3) , where the ordering of “ $* > A_1 > A_2 > A_3$ ” forms a high-to-low
 414 hierarchy, and so on. The minimal interesting layer (the m -layer) is (A_2, B_2, C_2) , and the
 415 o -layer is $(A_1, *, C_1)$. From the m -layer (the bottom cuboid) to the o -layer (the top-cuboid
 416 to be computed), there are in total $2 \times 3 \times 2 = 12$ cuboids, as shown in figure 3.

417 Suppose that the popular drilling path is given (which can usually be derived based on
 418 domain expert knowledge, query history, and statistical analysis of the sizes of intermediate
 419 cuboids). Assume that the given popular path is $\langle (A_1, *, C_1) \rightarrow (A_1, *, C_2) \rightarrow (A_2, *, C_2)$
 420 $\rightarrow (A_2, B_1, C_2) \rightarrow (A_2, B_2, C_2) \rangle$, shown as the darkened path in figure 3. Then each path of
 421 an H-tree from root to leaf is ordered the same as the popular path.

422 This ordering generates a compact tree because the set of low level nodes that share the
 423 same set of high level ancestors will share the same prefix path using the tree structure.
 424 Each tuple, which represents the currently in-flow stream data, after being generalized to
 425 the m -layer, is inserted into the corresponding path of the H-tree. An example H-tree is
 426 shown in figure 4. In the leaf node of each path, we store relevant measure information of

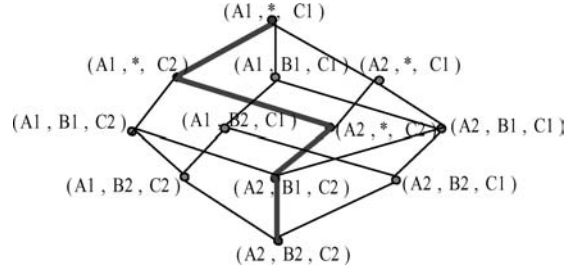


Figure 3. Cube structure from the m -layer to the o -layer.

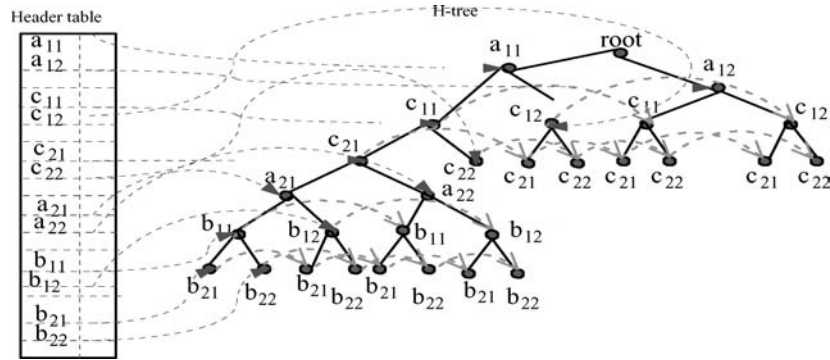


Figure 4. H-tree structure for cube computation.

the cells of the m -layer. The measures of the cells at the upper layers are computed using the H-tree and its associated links.

An obvious advantage of the *popular path approach* is that the nonleaf nodes represent the cells of those layers (cuboids) along the popular path. Thus these nonleaf nodes naturally serve as the cells of the cuboids along the path. That is, it serves as a data structure for intermediate computation as well as the storage area for the computed measures of the layers (i.e., cuboids) along the path.

Furthermore, the H-tree structure facilitates the computation of other cuboids or cells in those cuboids. When a query or a drill-down clicking requests to compute cells outside the popular path, one can find the closest lower level computed cells and use such intermediate computation results to compute the measures requested, because the corresponding cells can be found via a linked list of all the corresponding nodes contributing to the cells.

4.2. Algorithms for cube measure computation

With popular path stream data cube design and the H-tree data structure, the popular-path-based stream data cubing can be partitioned into three stages: (1) the initial computation of (partially materialized) stream data cube by popular-path approach, (2) incremental update

of stream data cube, and (3) online query answering with the popular-path-based stream data cube.

Here we present the three corresponding algorithms, one for each stage of the popular-path-based stream data cubing.

First, we present an algorithm for computation of initial (partially materialized) stream data cube by popular-path approach.

Algorithm 1 (Popular-path-based stream cube computation). Computing initial stream cube, i.e., the cuboids along the *popular-path* between the *m*-layer and the *o*-layer, based on the currently collected set of input stream data.

Input. (1) multi-dimensional multi-level stream data (which consists of a set of tuples, each carrying the corresponding time stamps), (2) the *m* and *o*-layer specifications, and (3) a given popular drilling path.

Output. All the aggregated cells of the cuboids along the popular path between the *m*- and *o*- layers.

Method.

1. Each tuple, which represents a minimal addressing unit of multi-dimensional multilevel stream data, is scanned once and generalized to the *m*-layer. The generalized tuple is then inserted into the corresponding path of the H-tree, increasing the count and aggregating the measure values of the corresponding leaf node in the corresponding slot of the tilted time frame.
2. Since each branch of the H-tree is organized in the same order as the specified popular path, aggregation for each corresponding slot in the tilted time frame is performed from the *m*-layer all the way up to the *o*-layer by aggregating along the popular path. The step-by-step aggregation is performed while inserting the new generalized tuples in the corresponding time slot.
3. The aggregated cells are stored in the nonleaf nodes in the H-tree, forming the computed cuboids along the popular path.

Analysis. The H-tree ordering is based on the popular drilling path given by users or experts. This ordering facilitates the computation and storage of the cuboids along the path. The aggregations along the drilling path from the *m*-layer to the *o*-layer are performed during the generalizing of the stream data to the *m*-layer, which takes only one scan of stream data. Since all the cells to be computed are the cuboids along the popular path, and the cuboids to be computed are the nonleaf nodes associated with the H-tree, both space and computation overheads are minimized.

Second, we discuss how to perform incremental update of the stream data cube in the popular-path cubing approach. Here we deal with the “always-grow” nature of time-series stream data in an “on-line,” continuously growing manner.

The process is essentially an incremental computation method illustrated below, using the tilted time frame of figure 1. Assuming that the memory contains the previously computed *m* and *o*-layers, plus the cuboids along the popular path, and stream data arrive every second. The new stream data are accumulated (by generalization) in the corresponding H-tree leaf nodes. If the time granularity of the *m*-layer is minute, at the end of every minute, the data

will be aggregated and be rolled up from leaf to the higher level cuboids. When reaching a cuboid whose time granularity is quarter, the rolled measure information remains in the corresponding minute slot until it reaches the full quarter (i.e., 15 minutes) and then it rolls up to even higher levels, and so on.

Notice in this process, the measure in the time interval of each cuboid will be accumulated and promoted to the corresponding coarser time granularity, when the accumulated data reaches the corresponding time boundary. For example, the measure information of every four quarters will be aggregated to one hour and be promoted to the hour slot, and in the mean time, the quarter slots will still retain sufficient information for quarter-based analysis. This design ensures that although the stream data flows in-and-out, measure always keeps up to the most recent granularity time unit at each layer.

We outline the incremental algorithm of the method as follows.

Algorithm 2 (Incremental update of popular-path stream cube with incoming stream data). Incremental computing stream cube, i.e., the cuboids along the *popular-path* between the *m*-layer and the *o*-layer, based on the previously computed cube and the newly input stream data.

Input. (1) a popular path-based stream data cube, which also includes (i) the *m* and *o*-layer specifications, and (ii) a given popular drilling path, and (2) a set of input multi-dimensional multi-level stream data (which consists of a set of tuples, each carrying the corresponding time stamps).

Output An updated stream data cube (i.e., the updated popular-path cuboids (between the *m*- and *o*-layers)).

Method.

1. Each newly coming tuple, which represents a minimal addressing unit of multi-dimensional multi-level stream data, is scanned once and generalized to the *m*-layer. The generalized tuple is then inserted into the corresponding path of the H-tree. If there exists a corresponding leaf node in the tree, increase the count and aggregating the measure values of the corresponding leaf node in the corresponding slot of the tilted time frame. If there exists no corresponding leaf node in the tree, a new leaf node is created in the corresponding path of the H-tree.
2. Since each branch of the H-tree is organized in the same order as the specified popular path, aggregation for each corresponding slot in the tilted time frame is performed from the *m*-layer all the way up to the *o*-layer by aggregating along the popular path. The step-by-step aggregation is performed while inserting the new generalized tuples finishes.
3. If it reaches the time when a sequence of data in the lower-level time slots should be aggregated to a new slot in the corresponding higher level titled time window, such aggregation will be performed at each level of the popular path. If it reaches the time when the data in the most distant time slot should be dropped from the valid time scope, the slot in the corresponding time window will be cleared.
4. The so computed aggregated cells are stored in the nonleaf nodes in the H-tree, forming the computed cuboids along the popular path.

527 Analysis. Based on our design of the tilted time window, such incremental computation can
 528 be performed along the popular path of the H-tree. Moreover, the aggregations along the
 529 drilling path from the m -layer to the o -layer are performed when the input stream data come
 530 to the m -layer, which takes only one scan of stream data. Since all the cells in the titled
 531 time windows in the cuboids along the popular path are incrementally updated, the cuboids
 532 so computed are correctly updated stream cube, with minimal space and computation
 533 over-heads.

534 Finally, we examine how fast online computation can be performed with such a partially
 535 materialized popular-path data cube. Since the query inquiring the information completely
 536 contained in the popular-path cuboids can be answered by directly retrieving the informa-
 537 tion stored in the popular-path cuboids, our discussion here will focus on retrieving the
 538 information involving the aggregate cells not contained in the popular-path cuboids.

539 A multi-dimensional multi-level stream query usually provides a few instantiated con-
 540 stants and inquires information related to one or a small number of dimensions. Thus one
 541 can consider a query involving a set of instantiated dimensions, $\{D_{ci}, \dots, D_{cj}\}$, and a set
 542 of inquired dimensions, $\{D_{qi}, \dots, D_{qk}\}$. The set of relevant dimensions, D_r , is the union
 543 of the sets of instantiated dimensions and the inquired dimensions. For maximal use of
 544 the precom-puted information available in the popular path cuboids, one needs to find the
 545 highest-level popular path cuboids that contains D_r . If one cannot find such a cuboid in the
 546 path, one will use the *base cuboid* at the m -layer to compute it. Then the computation can be
 547 performed by fetching the relevant data set from the so found cuboid and then computing
 548 the cuboid consisting of the inquired dimensions.

549 The online OLAP stream query processing algorithm is presented as follows.

550 **Algorithm 3 (Online processing of stream OLAP query).** Online processing of stream
 551 OLAP query given the precomputed stream data cube, i.e., the cuboids along the *popular-*
 552 *path* between the m -layer and the o -layer.

553 Input. (1) a popular path-based stream data cube, which includes (i) the m and o -layer
 554 specifications, and (ii) a given popular drilling path, and (2) a given query whose relevant
 555 dimension set is D_r , which in turn consists of a set of instantiated dimensions, $\{D_{ci}, \dots,$
 556 $D_{cj}\}$, and a set of inquired dimensions, $\{D_{qi}, \dots, D_{qk}\}$.

557 Output. A computed cuboid related to the stream OLAP query.

558 Method.

- 559 1. Find the highest-level popular path cuboids that contains D_r . If one cannot find such a
 560 cuboid in the path, one will use the *base cuboid* at the m -layer to compute it. Let the
 561 found cuboid be S .
- 562 2. Perform selection on S using the set of instantiated dimensions as set of constants,
 563 and using the set of inquired dimensions as projected attributed. Let S_c be the set of
 564 multidimensional data so selected.
- 565 3. Perform on line cubing on S_c and return the result.

566 Analysis. Based on our design of the stream data cube, the highest-level popular path
 567 cuboid that contains D_r should contain the answers we want. Using the set of instantiated
 568 dimensions as set of constants, and using the set of inquired dimensions as projected

569 attributed, the so-obtained S_c is the minimal set of aggregated data set for answering the
 570 query. Thus online cubing on this set of data will derive the correct result. Obviously, such
 571 a computation process makes good use of the precomputed cuboids and will involve small
 572 space and computation overheads.

573 5. Performance study

574 To evaluate the effectiveness and efficiency of our proposed stream cube and OLAP com-
 575 putation methods, we performed an extensive performance study on synthetic datasets. Our
 576 result shows that the total memory and computation time taken by the proposed algorithms
 577 are small, in comparison with several other alternatives, and it is realistic to compute such
 578 a partially aggregated cube, incrementally update them, and perform fast OLAP analysis
 579 of stream data using such precomputed cube.

580 Here we report our performance studies with synthetic data streams of various
 581 characteristics.² The data stream is generated by a data generator similar in spirit to the
 582 IBM data generator [5] designed for testing data mining algorithms. The convention for
 583 the data sets is as follows: *D3L3C10T400K* means there are 3 dimensions, each dimen-
 584 sion contains 3 levels (from the *m*-layer to the *o*-layer, inclusive), the node fan-out factor
 585 (cardinality) is 10 (i.e., 10 children per node), and there are in total 400 K merged *m*-layer
 586 tuples.

587 Notice that all the experiments are conducted in a static environment as a simulation of
 588 the online stream processing. This is because the cube computation, especially for full cube
 589 and top-k cube, may take much more time than the stream flow allows. If this is performed
 590 in the online streaming environment, substantial amount of stream data could have been
 591 lost due to the slow computation of such data cubes. This simulation serves our purpose
 592 since it clearly demonstrates the cost and the possible delays of stream cubing and indicates
 593 what could be the realistic choice if they were put in a dynamic streaming environment.

594 All experiments were conducted on a 2 GHz Pentium PC with 1 GB main memory,
 595 running Microsoft Windows-XP Server. All the methods were implemented using Sun
 596 Microsystems' Java 2 Platform, Standard Edition, version 1.4.2.

597 Our design framework has some obvious performance advantages over some alternatives
 598 in a few aspects, including (1) *tilted time frame* vs. *full non-tilted time frame*, (2) *using*
 599 *minimal interesting layer* vs. *examining stream data at the raw data layer*, and (3) *computing*
 600 *the cube up to the apex layer* vs. *computing it up to the observation layer*. Consequently,
 601 our feasibility study will not compare the design that does not have such advantages since
 602 they will be obvious losers.

603 Since a data analyst needs fast on-line response, and both space and time are critical in
 604 processing, we examine both time and space consumption. In our study, besides presenting
 605 the total time and memory taken to compute and store such a stream cube, we compare the
 606 two measures (time and space) of the *popular path* approach against two alternatives: (1)
 607 the *full-cubing* approach, i.e., materializing all the cuboids between the *m*- and *o*-layers,
 608 and (2) the *top-k cubing* approach, i.e., materializing only the top-k measured cells of the
 609 cuboids between the *m*- and *o*-layers, and we set top-k threshold to be 10%, i.e., only top
 610 10% (in measure) cells will be stored at each layer (cuboid). Notice that top-k cubing cannot

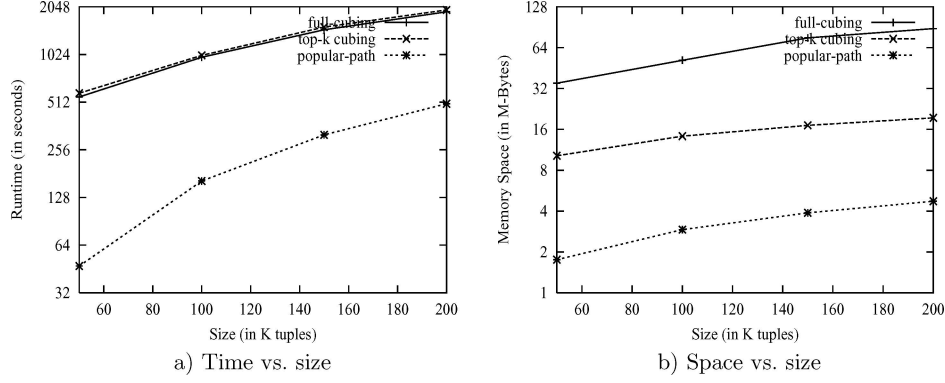


Figure 5. Cube computation: time and memory usage vs. no. tuples at the m -layer for the data set $D5L3C10$.

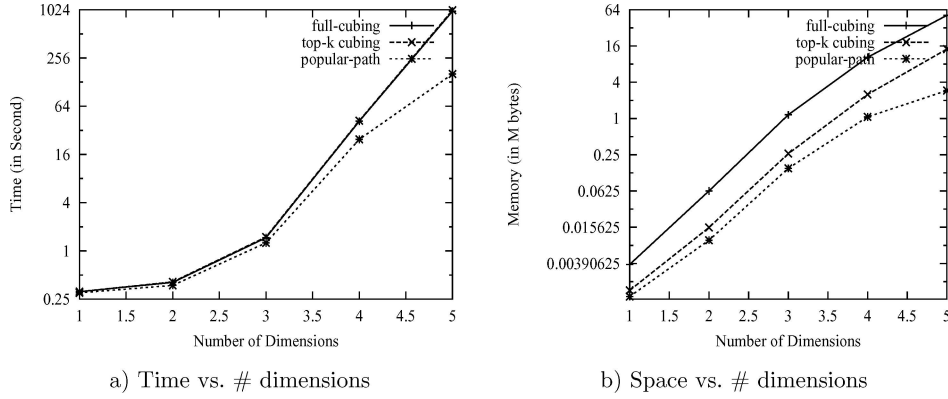


Figure 6. Cube computation: Time and space vs. no. of dimensions for the data set $L3C10/100K$.

be used for incremental stream cubing. However, since people may like to pay attention only to top-k cubes, we still put it into our performance study (as initial cube computation). From the performance results, one can see that if top-k cubing cannot compete with the popular path approach, with its difficulty at handling incremental updating, it will not likely be a choice for stream cubing architecture.

The performance results of stream data cubing (cube computation) are reported from figures 5 to 7.

Figure 5 shows the processing time and memory usage for the three approaches, with increasing size of the data set, where the size is measured as the number of tuples at the m -layer for the data set $D5L3C10$. Since *full-cubing* and *top-k cubing* compute all the cells from the m -layer all the way up to the o -layer, their total processing time is much higher than popular-path. Also, since *full-cubing* saves all the cube cells, its space consumption is much higher than popular-path. The memory usage of *top-k cubing* falls in between of the two approaches, and the concrete amount will depend on the k value.

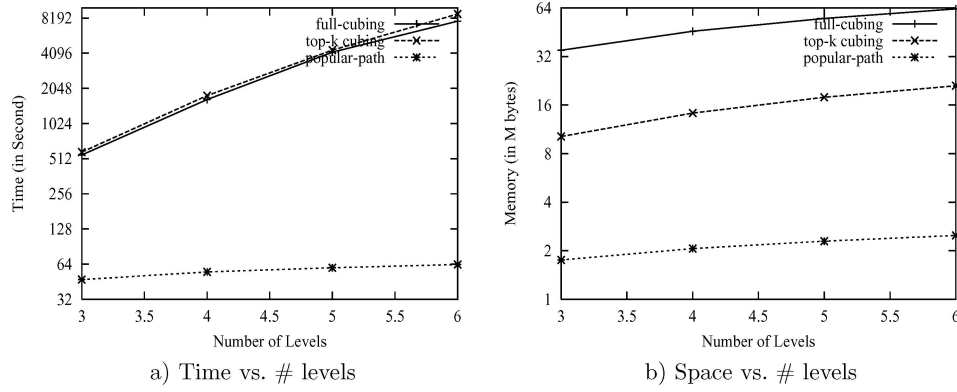


Figure 7. Cube computation: Time and space vs. no. of levels for the data set *D5C10T50K*. (a) Time vs. no. levels. (b) Space vs. no. levels.

Figure 6 shows the processing time and memory usage for the three approaches, with an increasing number of dimensions, for the data set *L3C10T100K*. Figure 7 shows the processing time and memory usage for the three approaches, with an increasing number of levels, for the data set *D5C10T50K*. The performance results show that popular-path is more efficient than both *full-cubing* and *top-k cubing* in computation time and memory usage. Moreover, one can see that increment of dimensions has a much stronger impact on the computation cost (both time and space) in comparison with the increment of levels.

Since incremental update of stream data cube carries the similar comparative costs for both popular-path and *full-cubing* approaches, and moreover, *top-k cubing* is inappropriate for incremental updating, we will not present this part of performance comparison. Notice that for incrementally computing the newly generated stream data, the computation time should be shorter than that shown here due to less number of cells involved in computation although the total memory usage may not reduce due to the need to store data in the layers along the popular path between two critical layers in the main memory.

Here we proceed to the performance study of stream query processing with four different approaches: (1) *full-cubing*, (2) *top-k cubing*, (3) popular-path, and (4) *no precomputation*, which computes the query and answer it on the fly. The reason that we added the fourth one is because one can compute query results without using any precomputed cube but using only the base cuboid: the set of merged tuples at the *m*-layer.

Figure 8 shows the processing time and memory usage vs. the size of the base cuboid, i.e., the number of merged tuples at the *m*-layer, for the data set *D5L3C10*, with the data set grows from 50 to 200 K tuples. There are 5 dimensions in the cube, and the query contains two instantiated columns and one inquired column. The performance results show that popular-path costs the least amount of time and space although *top-k cubing* could be a close rival. Moreover, *no precomputation*, though more costly than the previous two, still costs less in both time and space than the fully materialized stream cube at query processing.

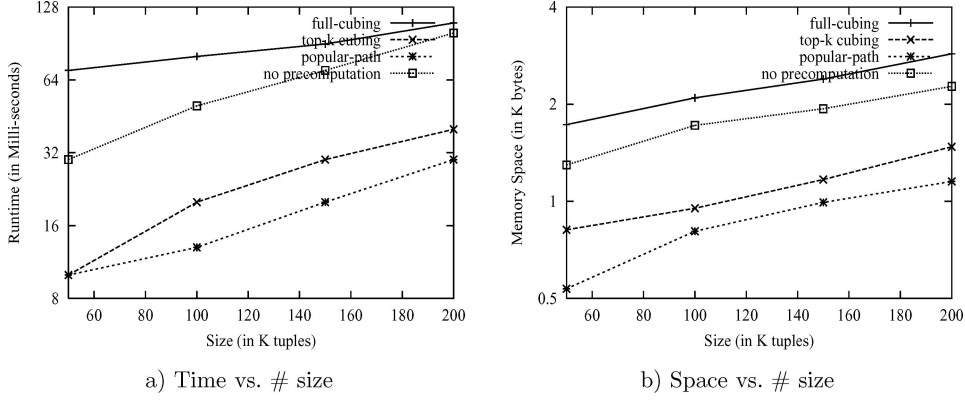


Figure 8. Stream query processing: Time and space vs. no. of tuples at the m -layer.

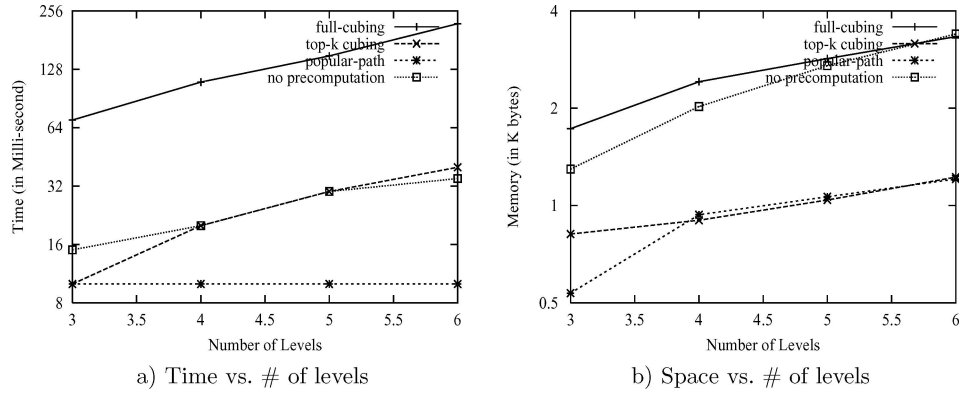


Figure 9. Stream query processing: Time and space vs. no. of levels.

Figure 9 shows the processing time and memory usage vs. the number of levels from the m to o layers, for the data set $D5C10T50K$, with the number of levels grows from 3 to 6. There are 5 dimensions in the cube, and the query contains two instantiated columns and one inquired column. The performance results show that popular-path costs the least amount of time and space and its query processing cost is almost irrelevant to the number of levels (but mainly relevant to the size of the tuples) with slightly increased memory usages. Moreover, *top-k cubing* and *no precomputation* takes more time and space when the number of levels increases. However, *full-cubing* takes the longest time to respond to a similar query although its response time is still in the order of 200 millisecond.

Finally, figure 10 shows the processing time and memory usage vs. the number of instantiated dimensions where the number of inquired dimensions maintains at one (i.e., single dimension) for the data set $D5L3C10T100K$. Notice that with more instantiated dimensions, the query processing cost for popular-path and *no precomputation* is actually

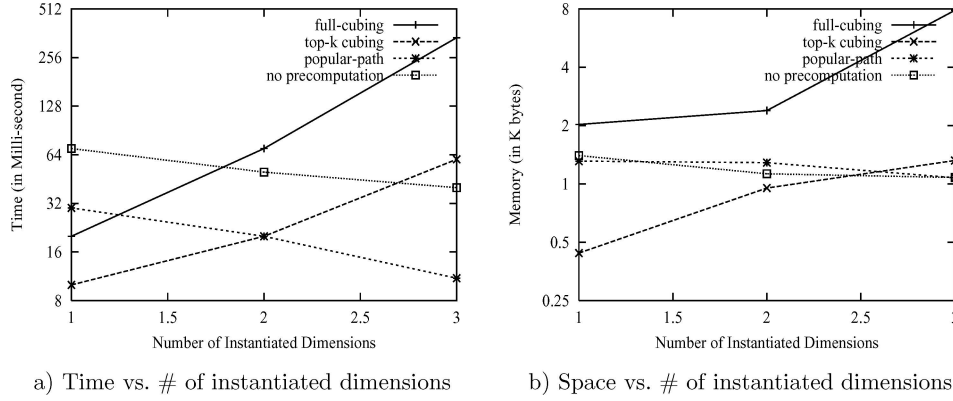


Figure 10. Stream query processing: Time and space vs. no. of instantiated dimension.

dropping because it will search less space in the H-tree or in the base cuboid with more instantiated constants. Initially (when the number of instantiated dimensions is only one, the *full-cubing* and *top-k cubing* are slightly faster than popular-path since the latter (popular-path) still needs some online computation while the former can fetch from the precomputed cubes.

From this study, one can see that popular-path is an efficient and feasible method for computing multi-dimensional, multi-level stream cubes, whereas *no precomputation* which computes only the base cuboid at the m -layer, could be the second choice. The full-cubing is too costly in both space and time, whereas *top-k cubing* is not a good candidate because it cannot handle incremental updating of a stream data cube.

6. Discussion

In this section, we compare our study with the related work and discuss some possible extensions.

6.1. Related work

Our work is related to: (1) on-line analytical processing and mining in data cubes, and (2) research into management and mining of stream data. We briefly review previous research in these areas and point out the differences from our work.

In data warehousing and OLAP, much progress has been made on the efficient support of standard and advanced OLAP queries in data cubes, including selective cube materialization [19], iceberg cubing [8, 18, 26, 28], cube gradient analysis [11, 21], exception [24], intelligent roll-up [25], and high-dimensional OLAP analysis [22]. However, previous studies do not consider the support for stream data, which needs to handle huge amount of fast changing stream data and restricts that the a data stream can be scanned only once. In contrast, our work considers complex measures in the form of stream data and studies OLAP and mining over partially materialized stream data cubes. Our data structure, to

certain extent, extend the previous work on H-tree and H-cubing [18]. However, instead of computing a materialized data cube as in H-cubing, we only use the H-tree structure to store a small number of cuboids along the popular path. This will save a substantial amount of computation time and storage space and lead to high performance in both cube computation and query processing. We have also studied whether it is appropriate to use other cube structures, such as star-trees in StarCubing [28], dense-sparse partitioning in MM-cubing [26] and shell-fragments in high-dimensional OLAP [22]. Our conclusion is that H-tree is still the most appropriate structure since most other structures need to either scan data sets more than once or know the sparse or dense parts beforehand, which does not fit the single-scan and dynamic nature of data streams.

Recently, there have been intensive studies on the management and querying of stream data [7, 12, 14, 16], and data mining (classification and clustering) on stream data [2–4, 13, 17, 20, 23, 27]. Although such studies lead to deep insight and interesting results on stream query processing and stream data mining, they do not address the issues of multidimensional, online analytical processing of stream data. Multidimensional stream data analysis is an essential step to understand the general statistics, trends and outliers as well as other data characteristics of online stream data and will play an essential role in stream data analysis. This study sets a framework and outlines an interesting approach to stream cubing and stream OLAP, and distinguishes itself from the previous work on stream query processing and stream data mining.

In general, we believe that this study sets a new direction: *extending data cube technology for multi-dimensional analysis of data streams*. This is a promising direction with many applications.

6.2. Possible extensions

There are many potential extensions of this work towards comprehensive, high performance analysis of data streams. Here we outline a few.

First, parallel and distributed processing can be used to extend the proposed algorithms in this promising direction to further enhance the processing power and the performance of the system. All of the three algorithms proposed in this study: *initial computation of stream data cubes*, *incremental update of stream data cube*, and *online multidimensional analysis of stream data*, can be handled by different processors and processed in a parallel and/or distributed manner. In the fast data streaming environment, it is desirable or sometimes required to have at least one processor dedicated to stream query processing (on the computed data cube) and at least another one dedicated to incremental update of data streams. Moreover, both incremental update and query processing can be processed by parallel processors as well since the algorithms can be easily transformed into parallel and/or distributed algorithms.

Second, although a stream cube usually retains in main memory for fast computation, updating, and accessing, it is important to have its important or substantial portion stored or mirrored on disk, which may enhance data reliability and system performance. There are several ways to do it. Based on the design of the tilted time frame, the distant time portion in the data cube can be stored on the disk. This may help reduce the total main

memory requirement and the update overhead. The incremental propagation of data in such distant portion can be done by other processors using other memory space. Alternatively, to ensure the data is not lost in case of system error or power failure, it is important to keep a mirror copy of the stream data cube on disk. Such a mirroring process can be processed in parallel by other processors. In addition, it is possible that a stream cube may miss a period of data due to software error, equipment malfunction, system failure, or other unexpected reasons. Thus a robust stream data cube should build the functionality to run despite the missing of a short period of data in the tilted time frame. The data so missed can be treated by special routines, like data smoothing, data cleaning, or other special handling so that the overall stream data can be interpreted correctly without interruption.

Third, although we did not discuss the computation of complex measures in the data cube environment, it is obvious that complex measures, such as sum, avg, min, max, last, standard deviation, linear regression and many other measures can be handled for the stream data cube in the same manner as discussed in this study. However, it is not clear how to handle holistic measures [15] in the stream data cubing environment. For example, it is still not clear how some holistic measures, such as quantiles, rank, median, and so on, can be computed efficiently in this framework. This issue is left for future research.

Fourth, the stream data that we discussed here are of simple numerical and categorical data types. In many applications, stream data may contain spatiotemporal and multimedia data. For example, monitoring moving vehicles and the flow of people in the airport may need to handle spatiotemporal and multimedia data. It is an open problem how to perform online analytical processing of multidimensional spatiotemporal and multimedia data in the context of data streams. We believe that spatiotemporal and multimedia analysis techniques should be integrated with our framework in order to make good progress in this direction.

Fifth, this study has been focused on multiple dimensional analysis of stream data. However, the framework so constructed, including tilted time dimension, monitoring the change of patterns in a large data cube using an m -layer and an o -layer, and paying special attentions on exception cells, is applicable to the analysis of non-stream time-series data as well.

Finally, this study is on multidimensional OLAP stream data analysis. Many data mining tasks require deeper analysis than simple OLAP analysis, such as classification, clustering and frequent pattern analysis. In principle, the general framework worked out in this study, including tilted time frame, minimal generalized layer and observation layers, as well as partial precomputation for powerful online analysis, will be useful for in-depth data mining methods. It is an interesting research theme on how to extend this framework towards online stream data mining.

7. Conclusions

In this paper, we have promoted on-line analytical processing of stream data, and proposed a feasible framework for on-line computation of multi-dimensional, multi-level stream cube.

We have proposed a general stream cube architecture and a stream data cubing method for on-line analysis of stream data. Our method uses a *tilted time frame*, explores *minimal*

775 *interesting and observation layers*, and adopts a *popular path approach* for efficient
 776 computation and storage of stream cube to facilitate OLAP analysis of stream data. Our
 777 performance study shows that the method is cost-efficient and is a realistic approach
 778 based on the current computer technology. Recently, this stream data cubing methodology
 779 has been successfully implemented in the MAIDS (Mining Alarming Incidents in
 780 Data Streams) project at NCSA (National Center for Supercomputing Applications) at
 781 University of Illinois, and its effectiveness has been tested using online stream data sets.

782 Our proposed stream cube architecture shows a promising direction for realization of
 783 on-line, multi-dimensional analysis of data streams. There are a lot of issues to be explored
 784 further. For example, besides H-cubing [18], there are other data cubing methodologies,
 785 such as multi-way array aggregation [29], BUC [8], and Star-cubing [28], it is interesting to
 786 examine other alternative methods for efficient online analysis of stream data. Furthermore,
 787 we believe that a very important direction is to further develop data mining methods to take
 788 advantage of multi-dimensional, multi-level stream cubes for single-scan on-line mining to
 789 discover knowledge in stream data.

790 Acknowledgments

791 The work was supported in part by research grants from U.S. National Science Foundation
 792 grants IIS-02-9199 and IIS-03-08215, Office of Naval Research, Natural Science and
 793 Engineering Research Council of Canada, and the University of Illinois. Any opinions,
 794 findings, and conclusions or recommendations expressed in this paper are those of the
 795 authors and do not necessarily reflect the views of the funding agencies. This paper is a
 796 substantially revised and major value-added version of a paper, “*Multi-Dimensional Re-*
 797 *gression Analysis of Time-Series Data Streams*,” by Yixin Chen, Guozhu Dong, Jiawei Han,
 798 Benjamin W. Wah, and Jianyong Wang, in VLDB’2002, Hong Kong, China, August, 2002.

799 Notes

- 800 1. We align the time axis with the natural calendar time. Thus, for each granularity level of the tilt time frame,
 801 there might be a partial interval which is less than a full unit at that level.
 802 2. We also tested it for some industry data sets and got similar performance results. However, we cannot discuss
 803 the results here due to the confidentiality of the data sets.

804 References

- 805 1. S. Agarwal, R. Agrawal, P.M. Deshpande, A. Gupta, J.F. Naughton, R. Ramakrishnan, and S. Sarawagi, “On
 806 the computation of multidimensional aggregates,” in Proc. 1996 Int. Conf. Very Large Data Bases (VLDB’96),
 807 Bombay, India, Sept. 1996, pp. 506–521.
 808 2. C. Aggarwal, J. Han, J. Wang, and P.S. Yu, “A framework for projected clustering of high dimensional data
 809 streams,” in Proc. 2004 Int. Conf. Very Large Data Bases (VLDB’04). Toronto, Canada, Aug. 2004, pp.
 810 852–863.
 811 3. C. Aggarwal, J. Han, J. Wang, and P.S. Yu, “On demand classification of data streams,” in Proc. 2004 Int.
 812 Conf. Knowledge Discovery and Data Mining (KDD’04), Seattle, WA, Aug. 2004, pp. 503–508.
 813 4. C.C. Aggarwal, J. Han, J. Wang, and P.S. Yu, “A framework for clustering evolving data streams,” in Proc.
 814 2003 Int. Conf. Very Large Data Bases (VLDB’03), Berlin, Germany, Sept. 2003.

- 815 5. R. Agrawal and R. Srikant, "Mining sequential patterns," in Proc. 1995 Int. Conf. Data Engineering (ICDE'95),
816 Taipei, Taiwan, March 1995, pp. 3–14.
- 817 6. B. Babcock, S. Babu, M. Datar, R. Motwani, and J. Widom, "Models and issues in data stream sys-
818 tems," in Proc. 2002 ACM Symp. Principles of Database Systems (PODS'02), Madison, WI, June 2002,
819 pp. 1–16.
- 820 7. S. Babu and J. Widom, "Continuous queries over data streams," SIGMOD Record, vol. 30, pp. 109–120,
821 2001.
- 822 8. K. Beyer and R. Ramakrishnan, "Bottom-up computation of sparse and iceberg cubes," in Proc. 1999 ACM-
823 SIGMOD Int. Conf. Management of Data (SIGMOD'99), Philadelphia, PA, June 1999, pp. 359–370.
- 824 9. S. Chaudhuri and U. Dayal, "An overview of data warehousing and OLAP technology," SIGMOD Record,
825 vol. 26, pp. 65–74, 1997.
- 826 10. Y. Chen, G. Dong, J. Han, B.W. Wah, and J. Wang, "Multi-dimensional regression analysis of time-series
827 data streams," in Proc. 2002 Int. Conf. Very Large Data Bases (VLDB'02), Hong Kong, China, Aug. 2002,
828 pp. 323–334.
- 829 11. G. Dong, J. Han, J. Lam, J. Pei, and K. Wang, "Mining multi-dimensional constrained gradients in data cubes,"
830 in Proc. 2001 Int. Conf. on Very Large Data Bases (VLDB'01), Rome, Italy, Sept. 2001, pp. 321–330.
- 831 12. J. Gehrke, F. Korn, and D. Srivastava, "On computing correlated aggregates over continuous data streams," in
832 Proc. 2001 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'01), Santa Barbara, CA, May 2001,
833 pp. 13–24.
- 834 13. C. Giannella, J. Han, J. Pei, X. Yan, and P.S. Yu, "Mining frequent patterns in data streams at multiple time
835 granularities," in H. Kargupta, A. Joshi, K. Sivakumar, and Y. Yesha (eds), Data Mining: Next Generation
836 Challenges and Future Directions. AAAI/MIT Press, 2004.
- 837 14. A.C. Gilbert, Y. Kotidis, S. Muthukrishnan, and M. Strauss, "Surfing wavelets on streams: One-pass summaries
838 for approximate aggregate queries," in Proc. 2001 Int. Conf. on Very Large Data Bases (VLDB'01), Rome,
839 Italy, Sept. 2001, pp. 79–88.
- 840 15. J. Gray, S. Chaudhuri, A. Bosworth, A. Layman, D. Reichart, M. Venkatrao, F. Pellow, and H. Pirahesh, "Data
841 cube: A relational aggregation operator generalizing group-by, cross-tab and sub-totals," Data Mining and
842 Knowledge Discovery, vol. 1, pp. 29–54, 1997.
- 843 16. M. Greenwald and S. Khanna, "Space-efficient online computation of quantile summaries," in Proc.
844 2001 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'01), Santa Barbara, CA, May 2001,
845 pp. 58–66.
- 846 17. S. Guha, N. Mishra, R. Motwani, and L. O'Callaghan, "Clustering data streams," in Proc. IEEE Symposium
847 on Foundations of Computer Science (FOCS'00), Redondo Beach, CA, 2000, pp. 359–366.
- 848 18. J. Han, J. Pei, G. Dong, and K. Wang, "Efficient computation of iceberg cubes with complex measures," in
849 Proc. 2001 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'01), Santa Barbara, CA, May 2001,
850 pp. 1–12.
- 851 19. V. Harinarayan, A. Rajaraman, and J.D. Ullman, "Implementing data cubes efficiently," in Proc. 1996 ACM-
852 SIGMOD Int. Conf. Management of Data (SIGMOD'96), Montreal, Canada, June 1996, pp. 205–216.
- 853 20. G. Hulten, L. Spencer, and P. Domingos, "Mining time-changing data streams," in Proc. 2001 ACM SIGKDD
854 Int. Conf. Knowledge Discovery in Databases (KDD'01), San Fransisco, CA, Aug. 2001.
- 855 21. T. Imielinski, L. Khachiyan, and A. Abdulghani, "Cubegrades: Generalizing association rules," Data Mining
856 and Knowledge Discovery, vol. 6, pp. 219–258, 2002.
- 857 22. X. Li, J. Han, and H. Gonzalez, "High-dimensional OLAP: A minimal cubing approach," in Proc. 2004 Int.
858 Conf. Very Large Data Bases (VLDB'04), Toronto, Canada, Aug. 2004, pp. 528–539.
- 859 23. G. Manku and R. Motwani, "Approximate frequency counts over data streams," in Proc. 2002 Int. Conf. Very
860 Large Data Bases (VLDB'02), Hong Kong, China, Aug. 2002, pp. 346–357.
- 861 24. S. Sarawagi, R. Agrawal, and N. Megiddo, "Discovery-driven exploration of OLAP data cubes," in Proc. Int.
862 Conf. of Extending Database Technology (EDBT'98), Valencia, Spain, March 1998, pp. 168–182.
- 863 25. G. Sathe and S. Sarawagi, "Intelligent rollups in multidimensional OLAP data," in Proc. 2001 Int. Conf. on
864 Very Large Data Bases (VLDB'01), Rome, Italy, Sept. 2001, pp. 531–540.
- 865 26. Z. Shao, J. Han, and D. Xin, "MM-Cubing: Computing iceberg cubes by factorizing the lattice space," in Proc.
866 2004 Int. Conf. on Scientific and Statistical Database Management (SSDBM'04), Santorini Island, Greece,
867 June 2004, pp. 213–222.

- 868** 27. H. Wang, W. Fan, P.S. Yu, and J. Han, "Mining concept-drifting data streams using ensemble classifiers," in
869 Proc. 2003 ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD'03), Washington, DC,
870 Aug. 2003.
- 871** 28. D. Xin, J. Han, X. Li, and B.W. Wah, "Star-cubing: Computing iceberg cubes by top-down and bottom-up
872 integration," in Proc. 2003 Int. Conf. Very Large Data Bases (VLDB'03), Berlin, Germany, Sept. 2003.
- 873** 29. Y. Zhao, P.M. Deshpande, and J.F. Naughton, "An array-based algorithm for simultaneous multidimensional
874 aggregates," in Proc. 1997 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'97), Tucson, Arizona,
875 May 1997, pp. 159–170.