# SMART: Scalable, Bandwidth-Aware Monitoring of Continuous Aggregation Queries

Navendu Jain, Praveen Yalagandula[†], Mike Dahlin, and Yin Zhang
University of Texas at Austin          [†]HP Labs

## ABSTRACT

We present SMART, a scalable, bandwidth-aware monitoring system that maximizes result precision of continuous aggregate queries over distributed data streams. While previous approaches *reduce bandwidth cost* under fixed precision constraints, in practice, monitoring systems may still incur a substantial cost risking overload under bursty traffic conditions. SMART therefore *bounds the worst-case system cost* to provide overload resilience and to facilitate practical deployment of monitoring systems. The primary challenge for SMART is how to select dynamic updates at each node in a distributed system to maximize global precision while keeping per-node monitoring bandwidth below a specified budget. To address this challenge, SMART's hierarchical algorithm (1) allocates bandwidth budgets in a near-optimal manner to maximize global precision and (2) self-tunes bandwidth settings to improve precision under dynamic workloads. Our prototype implementation of SMART provides key solutions to (a) prioritize pending updates for multi-attribute queries, (b) build bounded fan-in, load-aware aggregation trees to improve accuracy and fast anomaly detection, and (c) combine temporal batching with arithmetic filtering to reduce load and to quantify result staleness. Finally, our evaluation using simulations and a network monitoring application shows that SMART improves accuracy by up to an order of magnitude compared to uniform bandwidth allocation and performs close to the optimal algorithm under modest bandwidth budgets.

## 1. INTRODUCTION

Distributed stream processing systems [1,19,32] must provide high performance and high fidelity for query processing as such systems grow in scale and complexity. In these systems, data streams are often bursty where input rates may unexpectedly increase over time [3,23]. Examples include network traffic monitoring, identifying distributed denial-of-service (DDoS) attacks on the Internet, financial stocks monitoring, web click stream analysis, and event-driven monitoring in sensor networks. Therefore, it is desirable for these systems to bound the monitoring load while still providing useful accuracy guarantees on the query results. Existing techniques [21,22,24,26,37,40] aim to address this problem by minimizing the monitoring cost while promising an *a priori* numeric error bound (e.g., ± 10%) on query precision.

Unfortunately, although these techniques effectively reduce load under fixed precision, they are unsuitable in dynamic, high-volume stream processing environments for three reasons.

(1) **Setting precision requires workload knowledge:** Choosing error bounds *a priori* is unintuitive when workloads are not known in advance or may change unpredictably over time. (Should the error be 10% or 30%?) Conversely, it may be easier to set the monitoring budget (e.g., a system administrator is willing to pay 0.1% of network bandwidth for monitoring).

(2) **Bad precision setting hurts performance:** A bad choice of the error bound may significantly degrade the quality of a query result [6] (e.g., when the error bound is too large) or incur a high communication and processing cost (e.g., when the error bound is too small).

(3) **Bursty traffic imposes unacceptable overheads:** Even with reasonable error bounds, in practice, monitoring systems may still incur a substantial cost risking overload under bursty and often unpredictable traffic conditions (see motivating example described below).

An intuitive solution to handle high load is to simply provision adequate resources in anticipation of the worst-case load. However, since over-provisioning is neither economically viable [36] nor scalable [25], it is important to provide protection mechanisms to bound the bandwidth cost.

**Motivating Example:** We present a simple example to illustrate the challenges for scalable monitoring under bursty workloads. We simulate a set of 10 data sources connected to a centralized monitor with incoming bandwidth limit of up to 5 messages per second. The input workload distribution is modeled based on the standard exponential distribution with a parameter $\lambda$, and upon each arrival, the value of attribute $a_i$ at data source $i$ is updated according to random walk model in which the value either increases or decreases by an amount sampled uniformally from [0.5, 1.5]. Figure 1 shows the load-error tradeoff for a single data source. As expected, on increasing the error budget, arithmetic filtering [21,26] quickly decreases the load as majority updates get filtered.

To quantify the monitoring cost, we use $\lambda$=10, a baseline error budget of 2 and its corresponding expected load of 0.5 (Figure 1), effectively setting the total *expected* cost for monitoring 10 data sources as 5 messages per second. Figure 2 shows the *induced* message load at the central monitor under fixed error of 2 per attribute. We observe that under peak data arrivals, the system incurs up to 4x higher cost to meet the error bounds over the expected cost of 5 messages per second. Thus, monitoring systems may induce high overload to bound precision under bursty workloads. Figure 3 shows the divergence between data source attributes and
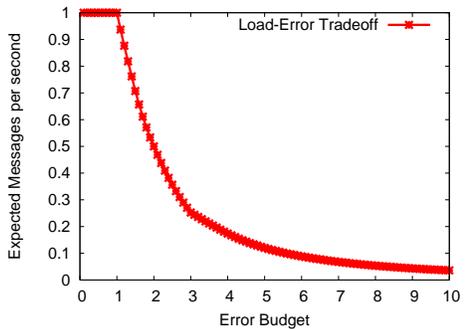
**Figure 1: Expected outgoing message load vs. error budget for a single attribute at a data source under random walk workload.**
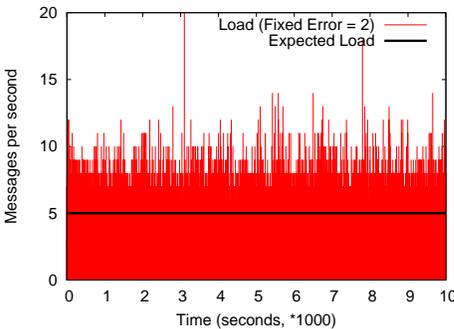
**Figure 2: Monitoring system induces overload under bursty workloads to bound result error. The system incurs up to 4x overload over the expected load.**
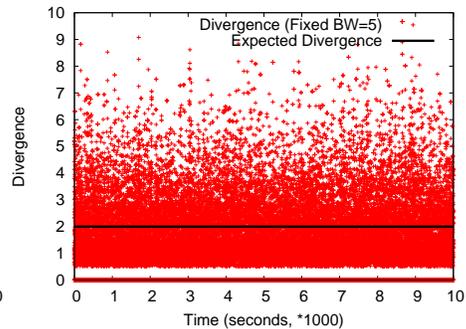
**Figure 3: Monitoring system causes high inaccuracy under bursty workloads to bound load. The result error is up to 5x higher over the expected divergence.**

their cached values at the central monitor under fixed load of 5. In this case, the system may provide highly inaccurate results having up to 5x error over the expected result error of 2. As a result, existing monitoring techniques are ineffective under bursty workloads—they risk overload, loss of accuracy, or both. This simple example clearly shows that large-scale monitoring systems must bound the bandwidth cost while still providing query results with useful accuracy guarantees.

**Our Contributions:** To address these challenges, we propose SMART, the first (to the best of our knowledge) scalable, bandwidth-aware monitoring system that adaptively sets bandwidth constraints to maximize precision of continuous aggregate queries over distributed data streams. SMART addresses the above challenges with the following techniques:

(1) **Maximize precision under fixed bandwidth budget:** SMART formulates a bandwidth-aware optimization problem whose goal is to maximize the result precision of an aggregate query in a hierarchical aggregation tree subject to given bandwidth constraints. This model provides a closed-form, near-optimal solution to self-tune bandwidth settings for achieving high accuracy using only local and aggregated information at each node in the tree. Further, SMART provides a "graceful degradation" in query precision as the available bandwidth decreases. Our experimental results show that self-tuning bandwidth settings improves accuracy by an order of magnitude over uniform bandwidth allocation and performs close to the optimal algorithm under modest bandwidth budgets.

(2) **Scalability via multiple aggregation trees:** SMART builds on recent work that uses distributed hash tables (DHTs) to construct scalable, load-balanced forests of self-organizing aggregation trees [5, 12, 29, 39]. Scalability to tens of thousands of nodes and millions of attributes is achieved by mapping different attributes to different trees. For each tree in this forest of aggregation trees, SMART's self-tuning algorithm adjusts bandwidth settings to achieve high precision under dynamic workloads where the estimate of load vs. precision trade-offs, and hence the optimal distribution, can change over time.

(3) **High performance by combining arithmetic filtering and temporal batching:** SMART integrates temporal batching with arithmetic filtering to reduce the monitoring load and to quantify staleness of query results. For high-volume workloads, it is advantageous for nodes to batch multiple updates (that arrive close in time) and send a single combined update, thereby more effectively using budgeted bandwidth for refreshing updates. In an aggregation tree, this temporal batching allows leaf sensors to condense a series of updates into a periodic report and allows internal nodes to combine updates from different subtrees before transmitting them further. Further, it bounds the delay (e.g., at most 30 seconds) from when an update occurs at a leaf until it is reported at the root. To minimize query result staleness, SMART prioritizes updates so that they propagate from the leaves to the root in the allotted time while meeting bandwidth budgets.

We have implemented a prototype of SMART on SDIMS, a scalable aggregation system [39] built on top of FreePastry [13]. Our prototype implementation provides two key optimizations. First, it constructs bounded fan-in, bandwidth-aware aggregation trees to improve accuracy and to quickly detect anomalies in heterogeneous environments. Recent studies [9, 11, 18, 21, 22] suggest that only a few attributes (e.g., network flows) generate a significant fraction of total traffic in many monitoring applications. Thus, for fast anomaly detection, an aggregation tree should quickly route their updates towards the root such that no internal node becomes a bottleneck due to either high in-degree or low bandwidth. Second, for multi-attribute queries, SMART provides a refresh schedule that selects and prioritizes attributes for refreshing in order to minimize the error in query results.

We evaluate SMART through extensive simulations and a real network monitoring application of detecting distributed heavy hitters. Experience with this application built on SMART illustrates the improved precision and scalability benefits: for a given monitoring budget, SMART's adaptivity can significantly improve the query precision while monitoring a large number of attributes. Compared to uniform bandwidth allocation, SMART improves accuracy by up to an order of magnitude and provides accuracy within 27% of the optimal algorithm under modest bandwidth budgets.

**Example Queries:** We list several real-world application queries which form the basis of SMART monitoring system.

$Q_1$  Monitor distributed heavy hitters that send the highest traffic in aggregate across all network endpoints.

$Q_2$  Find top-k ports across all nodes that have been heavily scanned in the recent past indicating worm activity.

$Q_3$  Monitor anomaly conditions e.g., SUM(nodes sensing fire) $\geq$ threshold, MAX(chemical concentration) in sensor networks.

$Q_4$  Monitor the top-k popular web objects in a wide-area content distribution network e.g., Akamai.

All these aggregate queries require processing a large number of rapid update streams in limited bandwidth/battery-life environments, and can benefit from SMART. In Section 5, we show results for $Q_1$ using a real network monitoring system we have implemented.

In summary, this paper makes the following contributions.

- The first research effort that identifies the key limitations of previous "fix error, minimize load" techniques and addresses them by bounding the worst-case load while still providing query results with useful accuracy.

- A practical, bandwidth-aware monitoring system that adapts bandwidth budgets to maximize precision of continuous aggregate queries under high-volume, dynamic workloads.

- Our implementation provides key optimizations for improving accuracy and fast anomaly detection in real-world heterogeneous environments.

- Our evaluation demonstrates that SMART provides a key substrate for scalable monitoring: it provides high accuracy in dynamic environments and performs close to the optimal algorithm under modest bandwidth budgets.

The rest of this paper is organized as follows. Section 2 provides background description of SDIMS [39], a scalable DHT-based aggregation system, and precision-performance tradeoffs that underlie SMART. Section 3 describes the SMART adaptive algorithm that self-tunes bandwidth settings to improve result accuracy. Section 4 presents the implementation of SMART in our SDIMS aggregation system. Section 5 presents the experimental evaluation of SMART. Finally, Section 6 discusses related work, and Section 7 provides conclusions.

## 2.  POINT OF DEPARTURE

SMART extends SDIMS [39] which embodies two key abstractions for scalable monitoring: aggregation and DHT-based aggregation. SMART then introduces controlled tradeoffs between precision bounds and monitoring load.

### 2.1  DHT-based Hierarchical Aggregation

Aggregation is a fundamental abstraction for scalable monitoring [5, 12, 18, 29, 37, 39] because it allows applications to access summary views of global information and detailed views of rare events and nearby information.
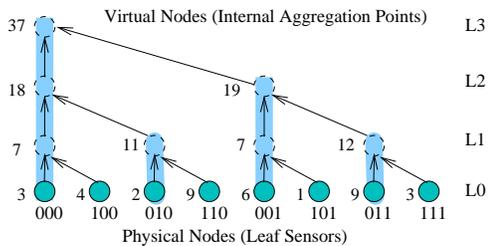


**Figure 4: The aggregation tree for key 000 in an eight node system. Also shown are the aggregate values for a simple SUM() aggregation function.**

SMART's aggregation abstraction defines a tree spanning all nodes in the system. As Figure 4 illustrates, each physical node is a leaf and each subtree represents a logical group of nodes[1]. An internal non-leaf node, which we call a *virtual node*, is simulated by a physical leaf node of the subtree rooted at the virtual node. Figure 4 illustrates the computation of a simple SUM aggregate.

SMART leverages DHTs [29,31,35] to construct a forest of aggregation trees and maps different attributes to different trees for scalability. DHT systems assign a long (e.g., 160 bits), random ID to each node and define a routing algorithm to send a request for ID $i$ to a node $root_i$ such that the union of paths from all nodes forms a tree $DHTtree_i$ rooted at the node $root_i$. By aggregating an attribute with ID $i$ = hash(attribute) along the aggregation tree corresponding to $DHTtree_i$, different attributes are load balanced across different trees. This approach can provide aggregation that scales to large numbers of nodes and attributes [5, 29, 39].

### 2.2  Query Result Approximation

SMART quantifies the precision of query results in terms of numeric error between the reported result and the actual value. We formally define the numeric approximation of a query result using Arithmetic Imprecision [21, 26].

*Arithmetic imprecision* (AI) deterministically bounds the numeric difference between a reported value of an attribute and its true value [21, 26, 27, 40]. For example, an AI = 10% bounds that the reported value either underestimates or overestimates the true value by at most 10%.

When applications do not need exact answers [21, 26, 37, 40], arithmetic imprecision provides an effective way to reduce system load by introducing additional filtering on update propagation. Next, we describe the AI mechanism of how SMART enforces the numeric error bounds while maximizing load reduction.

**AI Mechanism:** We first describe the basic mechanism for enforcing AI for each aggregation subtree in the system.

To enforce AI, each aggregation subtree $T$ for an attribute has an error budget $\delta_T$ that defines the maximum inaccuracy of any result the subtree will report to its parent for that attribute. The root of each subtree divides this error budget among itself $\delta_{self}$ and its children $\delta_c$ (with $\delta_T \geq \delta_{self} + \sum_{c \in children} \delta_c$), and the children recursively

---

[1]Logical groups can correspond to administrative domains (e.g., department or university) or groups of nodes within a domain (e.g., a /28 subnet with 14 hosts in the CS department) [16, 39].

do the same. Note that $\delta_c$ determines the child's error filter to cull as many updates as possible before sending to the parent, and $\delta_{self}$ is useful for applying additional filtering after combining all updates received from the children. Here we present the AI mechanism for the SUM aggregate since it is likely to be common in network monitoring [21, 26] and financial applications [15]; other standard aggregation functions (e.g., MAX, MIN, AVG, etc.) are similar and defined precisely in an extended technical report [20].

This arrangement reduces system load by filtering small updates that fall within the range of values cached by a subtree's parent. In particular, after a node A with error budget $\delta_T$ reports a range $[V_{min}, V_{max}]$ for an attribute value to its parent (where $V_{max} \leq V_{min} + \delta_T$), if the node A receives an update from a child, the node A can skip updating its parent as long as it can ensure that the true value of the attribute for the subtree lies between $V_{min}$ and $V_{max}$, i.e., if

$$
\begin{aligned}
V_{min} &\leq \sum_{c \in children} V^c_{min} \\
V_{max} &\geq \sum_{c \in children} V^c_{max}
\end{aligned}
\tag{1}
$$

where $V^c_{min}$ and $V^c_{max}$ denote the most recent update received from child $c$.

SMART maintains per-attribute $\delta_T$ values so that different attributes with different error requirements and different update patterns can use different $\delta$ budgets in different subtrees.

## 2.3 Case-study Application

To guide the system development of SMART and to drive our performance evaluation, we have built a case-study application using SDIMS aggregation framework: a distributed heavy hitter detection service. Distributed Heavy Hitters (DHH) detection is important for both monitoring traffic anomalies such as DDoS attacks, botnet attacks, and flash crowds as well as accounting and bandwidth provisioning [11].

Heavy hitters are entities that account for at least a specified proportion of the total activity measured in terms of number of packets, bytes, connections, etc. [11] in a distributed system—for example, the top 100 IPs that account for a significant fraction of total incoming traffic in the last 10 minutes [11]. To answer this distributed query, the key challenge is scalability for aggregating per-flow statistics for tens of thousands to millions of concurrent flows across all the network endpoints in real-time. For example, a subset of the Abilene [2] traces used in our experiments include 80 thousand flows that send about 25 million updates per hour. To process this workload, a centralized system needs to handle about 7000 messages per second at the central monitor.

To scalably compute the global heavy hitters list, we chain two aggregations where the results from the first feed into the second. First, SDIMS calculates the total incoming traffic for each destination address from all nodes in the system using SUM as the aggregation function and hash(HH-Step1, destIP) as the key. For example, tuple (H = hash(HH-Step1, 72.179.58.7), 900 KB) at the root of the aggregation tree $T_H$ indicates that a total of 900 KB of data was received for 72.179.58.7 across all vantage points in the network during the last time window. In the second step, we feed these aggregated total bandwidths for each destination IP into a SELECT-TOP-100 aggregation with key hash(HH-Step2, TOP-100) to identify the TOP-100 heavy hitters among all flows.

In Section 5 we show how SMART's self-tuning algorithm adapts bandwidth settings to monitor a large number of attributes and provides high result accuracy by filtering majority of mice flows [11] (attributes with low frequency) while prioritizing updates for the heavy-hitter flows; we expect typical monitoring applications to have a large number of mice flows but only a few heavy hitters.

## 3. SMART DESIGN

In this section we present the SMART design and describe our self-tuning algorithm that adapts bandwidth settings at each node to maximize query precision under given bandwidth budgets.

## 3.1 System Model

We focus on distributed stream processing environments with a large number of data sources that perform in-network aggregation to compute continuous aggregate queries over incoming data streams. The bandwidth resources for query processing may be limited at a number of points in the network. In particular, a node $j$'s *outgoing bandwidth* $(B^O_j)$ may be constrained, a node $j$'s *incoming bandwidth* $(B^I_j)$ may be constrained, or both. Note that constraining the incoming bandwidth bounds (a) the control traffic overhead for monitoring and (b) the CPU processing load for computing the aggregation function across incoming data inputs. As discussed in Section 1, bounding this incoming load is important for handling abnormal traffic conditions. Finally, these bandwidth capacities may vary (a) among nodes in heterogeneous environments and (b) with time if traffic is shared with other applications.

At any time, each SMART attribute's numeric value is bounded by an AI error $\delta$. If $\delta$ is small, then updates may frequently drive an attributes value out of its last reported range $[V_{min}, V_{max}]$, forcing the system to send messages to update the range. A system can, however, reduce its bandwidth requirements by increasing $\delta$. Thus, if a hierarchical aggregation system has bandwidth constraints, it should determine a $\delta$ value at each aggregation point that meets the bandwidth constraints into and out of that point, and it should select these $\delta$ values so as to minimize the total AI for the attribute. In particular, rather than splitting each node's incoming bandwidth evenly among its children, the system attempts to assign bandwidth to where it will do the most good by reducing the resulting imprecision of the node's aggregate values.

In the rest of this section, we describe the SMART algorithm for minimizing imprecision while meeting bandwidth constraints in four steps.

First, we describe a simplified algorithm for a one-level aggregation tree and static workloads. For each leaf node $i$ in the system, this algorithm calculates an ideal error setting $\delta_i$ and corresponding expected bandwidth consumption $b_i$ such that (1) each leaf node's outbound bandwidth (i.e., rate of updates sent to the root) is at most its outgoing bandwidth budget, (2) the root node's incoming bandwidth (i.e., sum of update rates inbound from children) is at most its incoming bandwidth budget, and (3) the sum of the $\delta_i$'s is minimized given the first two constraints.

Second, we describe how to handle dynamic workloads where the estimate of AI error $\delta$ vs. bandwidth trade-offs, and hence the optimal distribution, can change over time. A key challenge here is throttling the rate at which the system

| Symbol | Meaning |
|--------|---------|
| $B_i^O$ | outgoing bandwidth constraint for node $i$ |
| $B_i^I$ | incoming bandwidth constraint for node $i$ |
| $u_i$ | input update rate at node $i$ |
| $\sigma_i$ | standard deviation of node $i$'s input workload |
| $child(i)$ | all children of node $i$ |
| $\delta_i$ | node $i$'s AI error setting |
| $b_i$ | node $i$'s outgoing bandwidth load |
| $\delta_i^{opt}$ | node $i$'s optimal AI error setting |
| $b_i^{opt}$ | node $i$'s optimal outgoing bandwidth load |

**Table 1: Summary of key notations.**

redistributes bandwidth budgets across nodes since such redistribution also incurs bandwidth costs.

Third, we generalize the algorithm to handle multi-level aggregation trees.

Finally, we discuss how our implementation copes with variability. In particular, SMART sets the per-node $\delta$s so that the *average* bandwidth meets a target. However, spikes of update load for an attribute or coincident updates for multiple attributes could cause *instantaneous* bandwidth to exceed the target. To avoid such instantaneous overload, SMART therefore prioritizes pending updates based on the impact they will have on their aggregate values and drains them to the network at the target rate.

## 3.2 One-Level Tree

**Quantify AI Precision vs. Bandwidth Tradeoff:** To estimate the optimal distribution of load budgets among different nodes, we need a simple way of quantifying the amount of query error reduction that can be achieved when a given bandwidth budget is used for AI filtering.

Intuitively, the filtering gain depends on the size of the error budget relative to the inherent variability in the underlying data distribution. Specifically, as illustrated in Figure 5, if the precision threshold $\delta_i$ at node $i$ is much smaller than the standard deviation $\sigma_i$ of the underlying data distribution, $\delta_i$ is unlikely to filter many data updates but still consume valuable bandwidth. Meanwhile, if $\delta_i$ is above $\sigma_i$, we would expect the load to decrease quickly as $\delta_i$ increases until the point where a large fraction of updates are filtered.

To quantify the tradeoff between load and precision, we draw from the basic formulation of Jain et al. [21] to develop a simple metric in SMART for capturing the tradeoff between load and error budget. Our metric utilizes Chebyshev's inequality which gives a bound on the probability of deviation of a given random variable from its mathematical expectation in terms of its variance. Let $X$ be a random variable with finite mathematical expectation $\mu$ and variance $\sigma^2$. Chebyshev's inequality states that for any $k \geq 0$,

$$Pr(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2} \qquad (2)$$

For AI filtering, the term $k\sigma$ represents the error budget $\delta_i$ for node $i$. Substituting for $k$ in Equation 2 gives:

$$Pr(|X - \mu| \geq \delta_i) \leq \frac{\sigma_i^2}{\delta_i^2} \qquad (3)$$

Intuitively, this equation implies that if $\delta_i \leq \sigma_i$ i.e., the error budget is smaller than the standard deviation (implying $k \leq 1$), then $\delta_i$ is unlikely to filter many data updates (Figure 5.)

In this case, Equation 3 provides only a weak bound on the message cost: the probability that each incoming update will trigger an outgoing message is upper bounded by 1.
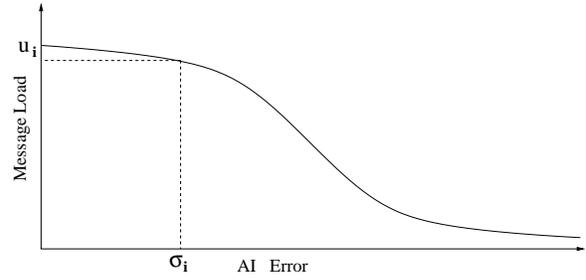


**Figure 5: Expected message load vs. AI error.**

However, if $\delta_i \geq k\sigma_i$ for any $k \geq 1$, the fraction of unfiltered updates is probabilistically bounded by $\frac{\sigma_i^2}{\delta_i^2}$. In general, given the input update rate $u_i$ for node $i$ with error budget $\delta_i$, the expected message cost for node $i$ per unit time is:

$$M_i = \text{MIN}\left\{u_i, \ \frac{\sigma_i^2}{\delta_i^2} * u_i\right\} \qquad (4)$$

We use the expected message cost $M_i$ to optimally set node $i$'s outgoing bandwidth $b_i^{opt}$.

**Estimate Optimal Bandwidth Settings under Fixed Load:** To estimate the optimal load distribution and settings of AI error $\delta$s at each node in a one-level tree rooted at node $r$, we formulate an optimization problem of minimizing the total error for a SUM aggregate at root $r$ under given bandwidth budgets $B_r^I$ (at the root) and $B_i^O$ (at child $i$). Specifically, we have:

$$\text{MIN} \sum_{i \in child(r)} \delta_i^{opt}$$
$$\text{s.t.} \sum_{i \in child(r)} \frac{\sigma_i^2 * u_i}{(\delta_i^{opt})^2} \leq B_r^I \qquad (5)$$
$$\forall i \in child(r), \quad b_i^{opt} = \frac{\sigma_i^2 * u_i}{(\delta_i^{opt})^2} \leq \min\{B_i^O, u_i\}$$

where $b_i^{opt}$ denotes the optimal setting of outgoing bandwidth of node $i$ to meet the global objective of minimizing the total error $\sum_{i \in child(r)} \delta_i^{opt}$ subject to the constraints of (1) node $i$'s outgoing bandwidth budget (i.e., $b_i^{opt} \leq \min\{B_i^O, u_i\}$) and (2) incoming bandwidth budget at root $r$ (i.e., $\sum_{i \in child(r)} b_i^{opt} \leq B_r^I$). This formulation is based on the AI filtering model i.e., node $i$ suppresses updates within $\delta_i$ of its last transmitted update.

**Enforcing Incoming Bandwidth Constraint:** To solve Equation (5), we first relax the outgoing bandwidth constraints $\forall i, \quad b_i^{opt} \leq \min\{B_i^O, u_i\}$ since the incoming bandwidth determines the processing overhead which seems likely to become a bigger bottleneck than the outgoing bandwidth. Later, we provide an optimal solution when the outgoing bandwidth constraints are also enforced. We use the method of Lagrangian Multipliers to find the extremum of the objective function $f : \left(\sum_i \delta_i^{opt}\right)$ subject to the constraint that

$$g : \left(\sum_i \frac{\sigma_i^2 * u_i}{(\delta_i^{opt})^2} - B_r^I\right) \leq 0.$$

The above formulation yields a closed-form and computa-

tionally inexpensive optimal solution [20]:

$$\delta_i^{opt} = \sqrt{\frac{\sum_{c\in child(r)} \sqrt[3]{\sigma_c^2 * u_c}}{B_r^I} * \sqrt[3]{\sigma_i^2 * u_i}} \qquad (6)$$

which provides a closed-form formula for setting $b_i^{opt}$:

$$b_i^{opt} = \frac{\sigma_i^2 * u_i}{(\delta_i^{opt})^2} = B_r^I * \frac{\sqrt[3]{\sigma_i^2 * u_i}}{\sum_{c\in child(r)} \sqrt[3]{\sigma_c^2 * u_c}} \qquad (7)$$

As a simple example, if $u_i = \sigma_i = 1$, then each node sets the same $\delta_i$ for the attribute as $\delta_i = \sqrt{\frac{N}{B_r^I}}$, and the bandwidth budget for node $i$ will be $\frac{1}{\delta_i^2} = \frac{B_r^I}{N}$ i.e., given a total bandwidth budget and uniform data distribution across nodes, each node gets a uniform share of the bandwidth to update each attribute.

Note that to set $b_i^{opt}$ (Equation (7)), each node needs to know $\sum_{c\in child(r)} \sqrt[3]{\sigma_c^2 * u_c}$ and root $r$'s incoming bandwidth budget $B_r^I$; SMART computes a simple SUM aggregate to obtain this information at each parent and propagates down to all its children in an aggregation tree. As a simple optimization, these messages are piggy-backed on data updates.

**Enforcing Outgoing Bandwidth Constraints:** Note that the above optimal load assignment assumes that the outgoing bandwidth constraint $b_c^{opt} \leq \min\{B_c^O, u_c\}$ holds for every child $c$. If for a node $i$, the above optimal solution doesn't satisfy $b_i^{opt} \leq \min\{B_i^O, u_i\}$, then we need to set $b_i^{opt} = \min\{B_i^O, u_i\}$ to satisfy its bandwidth constraint. This situation may arise in heterogeneous environments where a subset of nodes may experience larger input loads (e.g., DDoS attacks) or may become severely resource-constrained (e.g., sensor networks with low power devices).

Note that setting $b_i^{opt} = \min\{B_i^O, u_i\}$ can free up part of $r$'s incoming capacity $B_r^I$, which can be reassigned to other children to increase their bandwidth budget thereby improving the overall accuracy. SMART therefore applies an iterative algorithm that in each iteration determines all *saturated* children $C_{sat}$ at each step, fixes their load budgets and error settings, and recomputes Equations (6), (7) for all the remaining children (assuming child set $C_{sat}$ is absent). A child $j$ is labeled saturated if $b_j^{opt} \geq B_j^O$ i.e., the available outgoing bandwidth is less than or equal to that required by the optimal solution. If child set $C_{sat}$ is empty, the procedure terminates giving the optimal bandwidth and AI error settings for each node; all saturated nodes set bandwidth equal to their outgoing load threshold. Note that for our DHT-based aggregation trees, the fan-in for a node is typically 16 (i.e., a 4-bit correction per hop) so the iterative algorithm runs in constant time (at most 16 times).

### 3.3 Self-Tuning Bandwidth Settings

The above optimal solution is derived assuming that $\sigma_i$ and $u_i$ are given and remain constant. In practice, $\sigma_i$ and $u_i$ may change over time depending on the workload characteristics. SMART therefore performs self-tuning to dynamically readjust the bandwidth settings to minimize error.

**Relaxation:** A self-tuning algorithm that adapts too rapidly may react inappropriately to transient situations. We thus apply exponential smoothing to adjust the bandwidth set-

tings i.e.,

$$b_i^{new} = \alpha * b_i^{opt} + (1 - \alpha) * b_i \qquad (8)$$

where $\alpha = 0.05$.

**Cost-Benefit Throttling:** Finally, each node needs to send messages to the root to minimize the query error. Therefore, we need to prioritize sending those messages that benefit precision the most. A naive refreshing algorithm that updates attributes in a round-robin fashion could easily spend network messages that do not reduce error while consuming valuable bandwidth resource. Limiting sending of non-useful updates is a particular concern for applications like DHH that monitor a large number of attributes, only a few of which are active enough to be worth optimizing.

To address this challenge, after computing the new error budgets, a node computes a *charge* metric for each attribute $a$, which estimates the reduction in error gained by refreshing $a$:

$$charge_a = (T_{curr} - T_{lastSent}^a) * D_a \qquad (9)$$

where $T_{curr}$ is the current time, $T_{lastSent}^a$ is the last time an attribute $a$'s update was sent, and $D_a$ denotes the deviation between $a$'s current value at that node and its last reported range $[V_{min}, V_{max}]$ to the parent. For example, if $a$'s AI error range cached at the parent is [1, 2] and a new update with value 11 arrives at a child node, we expand the range to [1,11] at the child to include the new value setting $D_a =10$.

Notice that an attribute's charge will be large if (i) there is a large error imbalance (i.e., $D_a$ is large), or (ii) there is a long-lasting imbalance (e.g., $T_{curr} - T_{lastSent}^a$ is large). Further, only using $D_a$ may hurt precision if the attribute has a repeated behavior of quickly diverging after the last refresh. Therefore, the latter term prioritizes attributes who are likely to again diverge slowly after being refreshed thereby giving a long-term precision benefit.

Since redistribution also consumes bandwidth budgets, we only send messages to readjust bandwidth settings $b_i^{new}$ when doing so is likely to reduce the time-averaged error for those attributes by at least a threshold $\tau$ (i.e., if $charge_a > \tau$). Further, to ensure the invariant that a node does not exceed its bandwidth budget while sending updates for multiple attributes, we present a simple technique to prioritize updates across multiple attributes in Section 3.5.

### 3.4 Multi-Level Trees

To scale to a large number of nodes, we extend our basic algorithm for a one-level tree to a distributed algorithm for a multi-level aggregation hierarchy. Note that for a one-level aggregation tree, leaf nodes use AI error $\delta$ to filter updates. In addition, to reduce the communication and processing load in an aggregation hierarchy, the internal nodes may also retain a $\delta_{self}$ to help prevent updates received from their children from being propagated further up the tree [10,21].

At any internal node in the aggregation tree, the self-tuning algorithm works in a similar manner as that in a one-level tree case: the internal node is a *local* root for each of its immediate children, and the bandwidth targets $B^I$ for the parent and $B_c^O$ for each child $c$ are input as constraints in the optimization problem formulation in Equation (5). The key difference is that for children who are internal nodes, we use their AI error $\delta_{self}$ in this optimization framework. To estimate the optimal bandwidth and AI error settings for each child, the parent node $p$ tracks its incoming update rate

i.e., the aggregate number of messages sent by all its children per time unit ($u_p$) and the standard deviation ($\sigma_p$) of updates received from its children. Note that $u_c, \sigma_c$ reports are accumulated by child $c$ until they can be piggy-backed on an update message to its parent.

Given this information, the parent node estimates for each child $c$, the optimal outgoing bandwidth settings $b_c^{opt}$ and the corresponding AI error settings $\delta_{c(self)}^{opt}$ that minimize the total error in the aggregate value computed across the children given bandwidth constraints. This procedure is performed at each parent in the aggregation tree.

## 3.5 Prioritizing Pending Updates

Finally, we describe how our implementation addresses an important and practical challenge of variability in bandwidth targets. In particular, we want to avoid situations where a node may exceed *instantaneous* bandwidth target either due to (1) spikes of update load for an attribute, (2) coincident updates for multiple attributes, or (3) sudden increase in available bandwidth, but may still meet its *average* bandwidth target. This scenario is undesirable since a node may flood its parent with updates that far exceed the parent's incoming capacity.

To address this problem, our SMART implementation provides a priority heap that stores all pending updates that need to be sent ordered by their priority. Using this heap data structure, SMART efficiently prioritizes pending updates based on the impact they will have on their aggregate values and drains them to the network at a target rate. Specifically, at each time step, a node keeps removing the maximum priority attribute from the heap and sending it to the corresponding parent until the node's instantaneous target bandwidth is reached. Note that to compare priorities across different attributes that have different value ranges across different queries, we normalize an attribute $a$'s refresh priority (Equation (9)) with respect to their standard deviation by dividing the deviation $D_a$ by the standard deviation $\sigma_a$ of that attribute. Further, to support queries and attributes (within a query) with different importance values, SMART provides a simple mechanism of weighting an attribute's refresh priority with its importance value.

## 4. SMART IMPLEMENTATION

In this section, we describe several important optimizations and key issues for implementing SMART in our prototype implementation. First, we present a key performance optimization of temporal batching of updates to reduce load and to quantify staleness of query results in large-scale monitoring systems. Second, we describe how to build bandwidth-aware aggregation trees to improve accuracy and fast anomaly detection in heterogeneous network environments. Then, we describe how SMART handles failures and reconfigurations. Finally, we discuss how to improve precision for different aggregates.

## 4.1 Temporal Imprecision

SMART integrates temporal imprecision with arithmetic filtering to provide staleness guarantees on query results and to reduce the monitoring load.

*Temporal Imprecision* (TI) bounds the delay from when an event/update occurs until it is reported [24, 26, 34, 40]. A temporal imprecision of *TI* (e.g., TI = 30 seconds) guarantees that every event that occurred *TI* or more seconds



(a) Send synchronized updates every TI $-\Delta$ seconds.

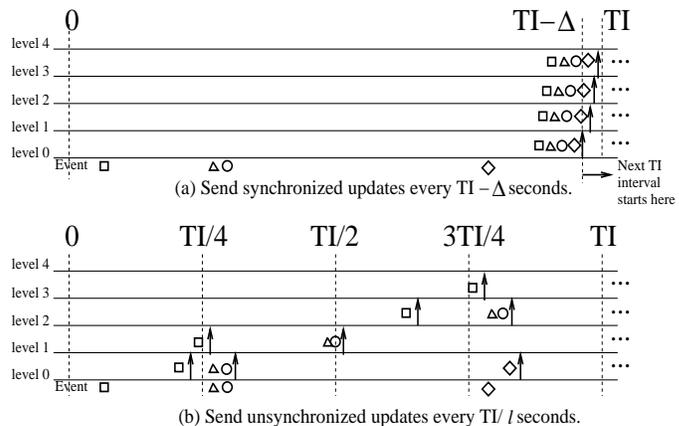(b) Send unsynchronized updates every TI/ *l* seconds.

**Figure 6: For a given TI bound, pipelined delays with synchronized clocks (a) allows nodes to send less frequently than unpipelined delays without synchronized clocks (b).**

ago is reflected in the reported result; events younger than *TI* may or may not be reflected. In SMART, each attribute has a TI interval during which its updates are batched into a combined message, checked if the combined update drives the aggregate value out of the last reported AI range, and then pushed into the priority queue to be sent to the parent.

Temporal imprecision benefits monitoring applications in two ways. First, it accounts for inherent network and processing delays in the system; given a worst case per-hop cost $hop_{max}$ even immediate propagation provides a temporal guarantee no better than $\ell * hop_{max}$ where $\ell$ is the maximum number of hops from any leaf to the root of an aggregation tree. Second, explicitly exposing TI allows SMART to reduce load by using temporal batching: a set of updates at a leaf sensor are condensed into a periodic report or a set of updates that arrive at an internal node over a time interval are combined into a single message before being sent further up the tree [20]. This temporal batching improves scalability by reducing processing and network load as we show using experiments on a network monitoring application in Section 5.

SMART implements TI using a simple mechanism of having each node send updates to its parent once per $TI/\ell$ seconds similar to TAG [24] as shown in Figure 6(b). Further, to maximize the possibility of batching updates, when clocks are synchronized[2], SMART pipelines delays across tree levels so that each node sends once every ($TI - \Delta$) seconds with each level's sending time staggered so that the updates from level $i$ arrive just before level $i + 1$ can send (Figure 6(a)). The term $\Delta$ accounts for the worst-case per hop delays and maximum clock skew; details are in the extended technical report [20].

## 4.2 Bandwidth-Aware Tree Construction

As described in Section 2, SMART leverages DHTs [29–31, 35] to construct a forest of aggregation trees and maps different attributes to different trees [5, 12, 29, 39] for scalability and load balancing. SMART then uses these trees to perform in-network aggregation.

---

[2] Algorithms in the literature can achieve clock synchronization among nodes to within one millisecond [38].

SMART constructs bounded fan-in, bandwidth-aware aggregation trees to improve result accuracy and quickly detect anomalies in heterogeneous environments. Recent studies [9, 11, 18, 21, 22] suggest that only a few attributes (e.g., elephant flows [11]) generate a significant fraction of total traffic in many monitoring applications. Thus, to provide fast anomaly detection, an aggregation tree should quickly route the updates of elephant flows towards the root such that no internal node becomes a processing/communication bottleneck due to either high in-degree or low bandwidth.

DHTs provide different degrees of flexibility in choosing neighbors and next-hop paths in building aggregation trees [14]. Many DHT implementations [13, 35] use proximity (usually round-trip latency) in the underlying network topology to select neighbors in the DHT overlay. This neighbor selection in turn determines the parent-child relationships in the aggregation tree. However, in a heterogeneous environment where different nodes have different bandwidth budgets, an aggregation tree formed solely based on RTTs may degrade the quality precision of the query result for two reasons. First, a node may not have sufficient outgoing bandwidth to send updates up in the tree even though its underlying tree may be sufficiently well-provisioned. In such an environment, this node becomes a bottleneck as the updates sent by the underlying subtree go wasted without benefiting the query precision. Second, a resource-limited parent may not be able to process the aggregate outgoing update rate of all its children. Thus, a practical protocol for building trees would be to bound the number of children at each internal node and use both latencies and bandwidth constraints to select the best parent at each tree level.

To improve accuracy and quickly identify anomalies, SMART builds DHT-based aggregation trees as follows:

- Bound the fan-in (i.e., number of children) at each parent node. In our implementation, a child node selects its parent such that the fan-in at the parent is at most 16 i.e., each parent has maximum up to 16 children.

- Use both network latency and available bandwidth capacity as the proximity metric for selecting parent nodes. SMART orders DHT neighbors of a node such that they have the highest incoming bandwidth capacity and have network latency below a specified threshold. Thus, nodes close in proximity and having high bandwidth capacities are highly likely to be selected as parent nodes.

Finally, note that in some environments, it might be useful to select nodes with low bandwidth as parents e.g., if the input workloads at the leaf nodes comes from an independent uniform distribution, then a node closer to the root is expected to receive very few updates since an aggregate (e.g., SUM) is likely to become more "stable" going towards the root. However, in practice, real workloads (1) are often non-uniform with few attributes generating a significant fraction of the total traffic [11] and (2) exhibit both temporal and spatial skewness with input rates unexpectedly increasing over time and across nodes. We quantify the effectiveness of constructing bounded fan-in, bandwidth-aware aggregation trees in Section 5.

## 4.3 Robustness

Failures and reconfigurations are common in large scale systems. As a result, a query might return a stale answer when nodes whose inputs are needed to compute the aggregate result become unreachable. More importantly, in a large scale monitoring system, such failures can interact badly with our techniques for providing scalability— hierarchy, arithmetic filtering, and temporal batching. For example, if a monitoring subtree is silent over an interval, it is difficult to distinguish between two cases: (a) the subtree has sent no updates because the inputs have not significantly changed or (b) the inputs have significantly changed but the subtree is unable to transmit its report. As a result, reported results can deviate arbitrarily from the truth.

Addressing this fundamental problem of node failures and network disruptions in large-scale monitoring systems is beyond the scope of this paper. In a separate study, we have developed a new consistency metric called Network Imprecision (NI) that characterizes and quantifies the accuracy of query results in the face of failures, network disruptions, and system reconfigurations; the details are available in an extended technical report [20].

## 4.4 Improving precision for different aggregates

SMART focuses on SUM aggregate since it is likely to be common in network monitoring [21, 26] and financial applications [15]; maximizing precision for the AVG aggregate is similar to SUM. For MAX, MIN aggregates, if the AI error budget is fixed, then the best error assignment is to give equal AI error budget to all the leaf nodes. However, since bandwidth budget is limited in practice, each node may set its precision differently to meet its available bandwidth. An intuitive solution to compute global MAX, MIN values is to simply broadcast them to each node, and a node sends an update only if it changes the global aggregate. However, this approach may limit scalability in large-scale data stream systems that monitor a large number of dynamic attributes. For the TOP-K aggregate, SMART achieves high accuracy by prioritizing updates based on the highest aggregate values and the result deviation $D_a$. Our experiments in Section 5 show that this approach works well in practice. We plan to develop mechanisms for other aggregates such as quantiles in future work.

## 5. EXPERIMENTAL EVALUATION

In this section, we present the precision and scalability results of an extensive experimental study of our SMART algorithm in a data streaming environment. First, we use simulations to evaluate the improvement in query result accuracy due to SMART's adaptive bandwidth settings. Second, we quantify the accuracy achieved by SMART for the DHH application in a network monitoring implementation. To perform this experiment, we implemented a prototype of SMART using SDIMS aggregation system [39] on top of FreePastry [31]. We used Abilene [2] netflow traces and performed the evaluation on 120 node instances mapped to 30 physical machines in the department Condor cluster. Finally, we investigate the precision benefits of constructing bandwidth-aware tree construction using our prototype.

In summary, our experimental results show that SMART is an important and effective substrate for scalable monitoring: SMART provides high result accuracy while bounding the monitoring load, continuously adapts to dynamic workloads, and achieves significant precision benefits for an important real-world monitoring application of detecting dis-
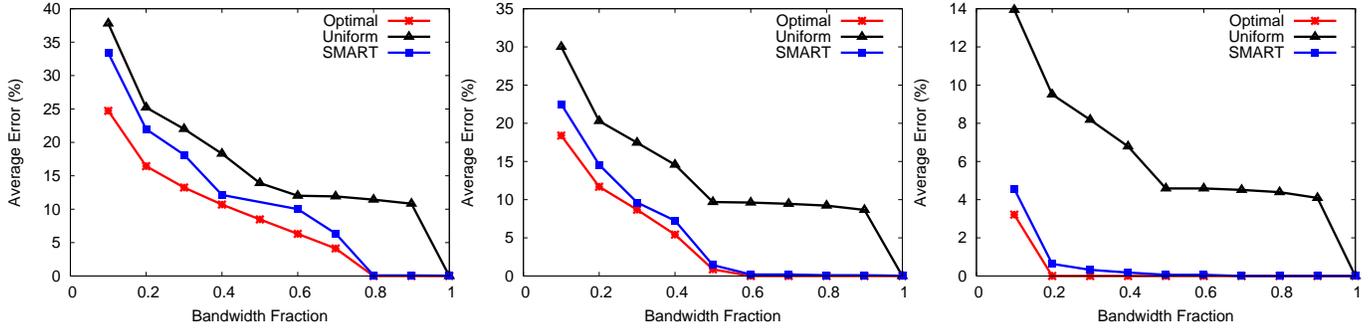
Figure 7: SMART provides higher precision benefits as skewness in a workload increases. The three figures show average error vs. load budget for three different skewness settings (a) 20:80% (b) 50:50% (c) 90:10%.

tributed heavy hitters.

## 5.1 Simulation Experiments

In this section, we quantify the result accuracy achieved by SMART compared to uniform bandwidth allocation and an idealized optimal algorithm. First, we assess the effectiveness of adaptive bandwidth settings in improving result precision as skewness in a workload increases. Second, we analyze the effect of fluctuating bandwidth. Finally, we evaluate SMART for different workloads.

In all experiments, all active sensors are at the leaf nodes of an aggregation tree. Each sensor generates a data value every time unit (round) for two sets of synthetic workloads for 100,000 rounds: (1) a random walk pattern in which the value either increases or decreases by an amount sampled uniformly from [0.5, 1.5], and (2) a Gaussian distribution with standard deviation 1 and mean 0. We simulate $m \in \{100, 1000\}$ data sources each having $n \in \{10, 100\}$ attributes in one-level and hierarchical topologies and under fixed and fluctuating bandwidth loads. All attributes have equal weights, messages have the same size, and each message uses one unit of bandwidth.

**Evaluating Update Rate Skewness:** First, we evaluate the precision benefits of SMART compared to other approaches as skewness in a workload increases. We compare it with (1) the optimal algorithm under the idealized and unrealistic model of perfect global knowledge of each attribute's divergence at each data source and (2) a uniform allocation policy where the incoming bandwidth capacity $B$ at a parent is allocated equally ($\frac{B}{c}$) among its $c$ children. Although this simple policy is correct (the total incoming load from the children is guaranteed to never exceed the incoming load of their parent at all times), it is not generally the best policy as we show in our experiments.

We first perform a simple experiment for a one-level tree, and later show the results for general hierarchical topologies. Figure 7 shows the query precision achieved by SMART for $m$=100 and $n$=10 under the following skewness settings: (a) 20:80% (b) 50:50% (c) 90:10%. For example, the 20:80% skewness represents that randomly selected 20% attributes are updated with probability 0.01 while the remaining ones are updated consistently every second under the random walk model. In all subsequent graphs in this section, the x-axis denotes the bandwidth budget as a fraction of the total cost $m.n$ of refreshing all the attributes across all nodes; the y-axis shows the resulting average error. For 20:80% skewness, since only a small fraction of attributes are stable,

SMART can only reclaim up to 20% load budget from stable attributes sources and distribute it to dynamic sources to reduce their error. For small bandwidth budgets, SMART improves accuracy by up to 35% compared to uniform allocation. The optimal algorithm improves accuracy by 27% over SMART. As the load budget increases, SMART converges to the optimal solution. SMART improves error by 40% over uniform allocation under 20% load budget and by more than an order of magnitude under sufficiently large budgets. For the 50:50 case, SMART can reclaim 50% of the total load budget compared to uniform allocation and give it to unstable sources. In this case, SMART reduces error by up to 50% over uniform policy at 40% bandwidth and achieves almost the accuracy of the optimal solution. Finally, for 90% skewness, SMART achieves the accuracy of the optimal algorithm even under 20% fraction of the total bandwidth and improves accuracy by more than an order of magnitude over uniform allocation. We observed qualitatively similar results for other $m$ and $n$ settings.

Note that the advantage of SMART's self-tuning algorithm depends on the skewness in the workload. We expect that for systems monitoring a large numbers of attributes (e.g., top-k heavy hitters query), some attributes (e.g., the elephants) will have high variability in data values and update rates so these attributes gain only a modest advantage in accuracy from SMART, while other attributes (e.g., the mice) will have large ratios and hence, the query will gain large advantages since a top-k query only needs to provide high accuracy for the top-k flows. We typically expect many more mice attributes than elephant attributes for common monitoring applications.

**Effect of Fluctuating Bandwidth:** Next, we evaluate the effectiveness of SMART under fluctuating bandwidth. We vary the incoming bandwidth over time following a sine wave pattern and set the maximum rate of bandwidth change to 10%, 20%, and 30% for $m = 1000$ nodes each having $n = 100$ attributes. We use update rate skewness of 50% as described above. From Figure 8, we observe that for 10% variation, SMART provides 50% reduction in error over uniform allocation at 40% bandwidth fraction. As we increase the bandwidth fluctuation from 10% to 30%, SMART reduces error by about 70% under 40% bandwidth fraction, and more than an order of magnitude for larger fractions. Further, in all cases, SMART achieves accuracy close to the optimal algorithm.

**Evaluating Different Workloads:** Finally, we evaluate the performance of SMART under different configurations
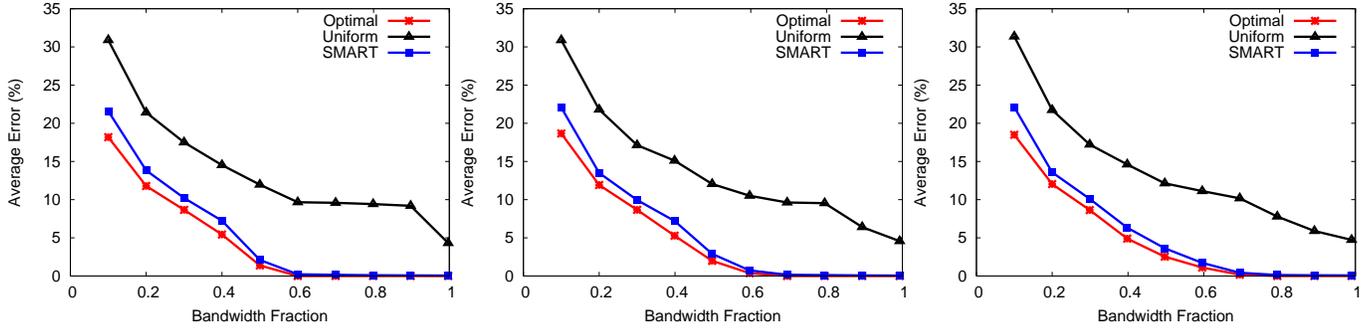
**Figure 8: SMART effectively reduces inaccuracy even under fluctuating bandwidth. The maximum bandwidth variation in the figure is (a) 10% (b) 20% and (c) 30%.**

by varying input data distribution, standard deviation (step sizes), and update frequency at each node. The workload data distribution is generated from a random walk pattern and Gaussian models. For standard deviation/step-size, 70% nodes have uniform parameters as previously described; the remaining 30% nodes have these parameters proportional to *rank* (i.e., with locality) or randomly assigned (i.e., no locality) from the range [0.5, 150].

Figure 9 shows the corresponding results for different settings of data distribution and standard deviation for $m=100$ and $n=100$. The update frequency is set to 0.7 skewness as described previously. We make three key observations. First, SMART minimizes error close to the optimum algorithm under the rank based assignment (Figure 9(a),(c)). Second, under random assignment, SMART achieves lesser accuracy benefits since updates generated from within the same subtree are not correlated. In both cases, as the bandwidth increases, SMART quickly minimizes error close to the optimal. Third, because step-sizes are based on node rank, SMART prioritizes attributes having the largest step-sizes and applies cost-benefit throttling to ensure that the precision benefits exceed costs. The uniform policy, however, does not make such a distinction equally favoring all attributes that need to be refreshed, thereby incurring a high error. Finally, under limited bandwidth, refreshing mice attributes with small step sizes does not significantly reduce the result error for queries such as TOP-K heavy hitters but consumes valuable bandwidth resources. For all these configurations, SMART reduces error by up to an order of magnitude over uniform allocation. The optimal approach reduces error by 20% over SMART.

Overall, across all configurations in Section 5.1, SMART reduces inaccuracy by up to an order of magnitude compared to uniform allocation and is within 27% of the optimal algorithm under modest bandwidth fraction.

## 5.2 Testbed Experiments

In this section, we quantify the error reduction of reported results due to self-tuning precision for the heavy hitter monitoring application.

We use multiple netflow traces obtained from the Abilene [2] Internet2 backbone network. The traces were collected from 3 Abilene routers for 1 hour; each router logged per-flow data every 5 minutes, and we split these logs into 120 buckets based on the hash of source IP. As described in Section 2.3, our DHH application executes a Top-100 heavy hitter query on this dataset for tracking the top 100 flows
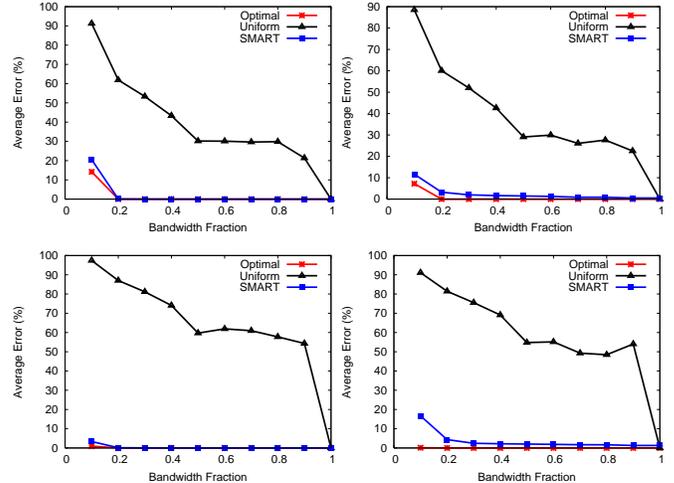


**Figure 9: Precision benefits of SMART vs. optimal algorithm and uniform allocation for different {workload, step sizes/standard deviation} configurations: (a) random walk, rank (b) random walk, random (c) Gaussian, rank, (d) Gaussian, random.**

(destination IP as key) in terms of bytes received over a 30 second moving window shifted every 10 seconds. We analyzed this workload [20] and observed that the 120 sensors track roughly 80,000 flows and send around 25 million updates in an hour. Further, it shows a heavy-tailed Zipf-like distribution: 60% flows send less than 1 KB of aggregate traffic, 90% flows less than 55 KB, and 99% of the flows send less than 330 KB during the 1-hour run; the maximum aggregate flow value is about 179.4 MB. We observed a similar heavy-tailed distribution for the number of updates per flow (attribute) [20].

For this experiment, we fixed the outgoing bandwidth as a constant between 0.5 and 10 messages per node per second. Since, we bound the fan-in of an internal node in our DHT-based aggregation tree to 16, the maximum incoming load at any node is thus 160 messages per second which is a reasonable processing load in our environment. Figure 10 plots the outgoing load per node on the x-axis and the result precision achieved for the top-100 heavy hitter query on the y-axis. The different lines in the graph correspond to a temporal batching interval of 10 seconds, 30 seconds, 60 seconds, and five minutes. Each data point denotes the average result divergence for the TOP-100 heavy hitters set.
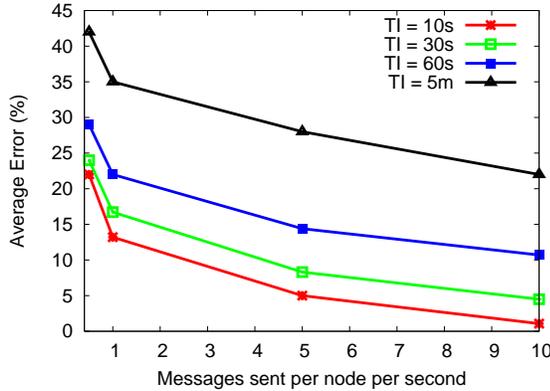
**Figure 10: Bandwidth load vs. query result precision for the DHH application.**
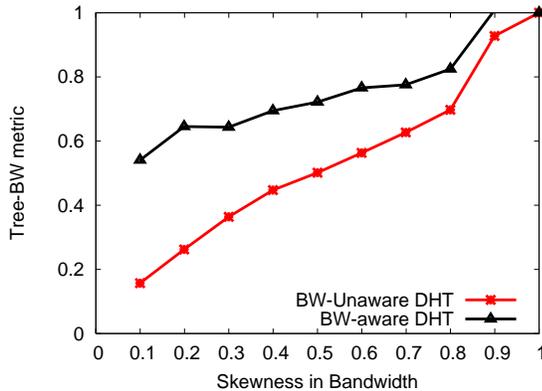


**Figure 11: Bandwidth skewness vs. normalized Tree-BW metric for bandwidth-aware and -unaware DHTs.**

It is important to note that for this query, the result set was consistent under TI of 10, 30, and 60 seconds with the true result set. However, under TI of 5 minutes, 5% results differed by less than 3%.

We observe that under TI = 10 seconds, SMART can reduce average error from 23% at 0.5 load budget to about 2% at load budget of 10. Comparing across TI values, the average error increases with increase in the TI batching interval. However, we believe that for many monitoring environments, the performance benefits of TI will outweigh the increased error as the bandwidth load can be significantly decreased due to temporal batching.

## 5.3 Evaluating Bandwidth-aware Tree Formation

In this section, we evaluate the benefits of SMART's bandwidth-aware tree construction. As described in Section 4, SMART uses both bandwidth and latency as proximity metrics in DHT routing compared to only latency in traditional DHT implementations. For quantitative comparison, we compute a *Tree-BW* metric for a tree as the sum of $L_i * B_i$ across all nodes, where $L_i$ is the number of leaves in the subtree rooted at node $i$ and $B_i$ is the bandwidth of node $i$. This weighted summation metric is higher for trees that select internal nodes having higher bandwidth capacities. For this experiment, we classify nodes belonging to two different

classes of bandwidth budgets: 100Mbps and 1Mbps. A 0.1 skewness in bandwidth implies that 10% of the nodes have 100Mbps bandwidth and the remaining have 1Mbps bandwidth. Figure 11 compares the benefits of SMART's tree construction using bandwidth-aware DHT routing against trees constructed in a bandwidth-oblivious unaware for a 1024-node system; the y-axis shows the normalized Tree-BW (with respect to the maximum value observed) metric for various bandwidth skewness settings (x-axis). Note that SMART's bandwidth-aware DHT routing achieves better Tree-BW metric values by up to a factor of 3.7x over a bandwidth-unaware DHT.

## 6. RELATED WORK

SMART is motivated by prior "fix error, minimize load" techniques [21, 22, 24, 26, 37, 40], but it departs in three significant ways driven by our focus on providing overload resilience for practical, scalable monitoring. First, SMART reformulates the key optimization problem, which we believe is an important contribution. While prior approaches "reduce cost" under fixed precision constraints, in practice, bursty, high-volume workloads still risk overloading a monitoring system as discussed in Section 1. Infact, it is precisely during these abnormal events that a monitoring system needs to provide high accuracy results with a real-time query response. SMART therefore "bounds the worst-case system cost" to provide overload resilience. Second, the technical advances to solve this new problem are significant. While SMART uses Chebyshev inequality [21] to capture the load vs. error trade-off, it faces new constraints of limited bandwidth budgets both at a parent and each of its children in a hierarchical aggregation tree. Given these constraints, SMART self-tunes bandwidth settings in a near-optimal manner to achieve high accuracy under dynamic workloads. Third, SMART's real-world implementation provides key solutions to improve result accuracy and fast anomaly detection in monitoring systems.

Recently, load shedding techniques [3, 36] have been proposed to handle system overload conditions. The key idea is to carefully drop some tuples to reduce processing load but at the expense of reducing the accuracy of query answers. Further, these approaches handle load spikes assuming either CPU [3, 36] or the main memory [17] as the key resource bottleneck. In comparison, bandwidth is the primary resource constraint in SMART.

Best-effort cache synchronization techniques [7,28] aim to minimize the divergence between source data objects and cached copies in a *one-level* tree. In these techniques, each object is treated individually for refreshing. In comparison, SMART performs hierarchical query processing for scalability and it proves valuable to co-relate updates to the same attribute at different data sources in order to maximize accuracy of an aggregate query result. To our knowledge, there has not been prior work on maximizing precision of aggregate queries in hierarchical trees under limited bandwidth.

Babcock and Olston [4] focus on efficiently computing top k aggregate values given fixed error in a single-level tree, but do not consider how to maximize precision of top k results under fixed bandwidth in hierarchical trees. Silberstein et al. propose a sampling-based approach (combined with global knowledge) at randomly chosen time steps for computing top-k queries in sensor networks [33]. Their approach focuses on returning the $k$ nodes with the highest

sensor readings In comparison, SMART focuses on large-scale aggregation queries such as distributed heavy hitters which require computing a total aggregate value for each of the tens of thousands to millions of attribute across all the nodes in the network. Finally, given the unpredictable traffic nature of network anomalies (e.g., DDoS attacks, flash crowds, botnet attacks), samples based on past readings are unlikely to be effective.

Sketches are small-space data structures that provide approximate answers to aggregate queries [8]. These techniques, however, require error bounds to be set a priori to provide the approximation guarantees.

# 7. CONCLUSIONS

We designed, implemented, and evaluated SMART—a scalable, load-aware monitoring system that performs self-tuning of bandwidth budgets to maximize precision of continuous aggregate queries.

In future work, we plan to examine techniques for secure information aggregation in distributed stream processing environments spanning multiple administrative domains as well as develop a broad range of distributed monitoring applications that could benefit from SMART.

# 8. REFERENCES

[1] D. J. Abadi, Y. Ahmad, M. Balazinska, U. Centintemel, M. Cherniack, J.-H. Hwang, W. Lindner, A. S. Maskey, A. Rasin, E. Ryvkina, N. Tatbul, Y. Xing, and S. Zdonik. The Design of the Borealis Stream Processing Engine. In *CIDR*, Asilomar, California, 2005.

[2] Abilene internet2 network. `http://abilene.internet2.edu/`.

[3] B. Babcock, M. Datar, and R. Motwani. Load shedding for aggregation queries over data streams. In *ICDE*, 2004.

[4] B. Babcock and C. Olston. Distributed top-k monitoring. In *SIGMOD*, June 2003.

[5] A. Bharambe, M. Agrawal, and S. Seshan. Mercury: Supporting Scalable Multi-Attribute Range Queries. In *SIGCOMM*, Portland, OR, August 2004.

[6] R. Cheng, B. Kao, S. Prabhakar, A. Kwan, and Y. Tu. Adaptive stream filters for entity-based queries with non-value tolerance. In *VLDB*, 2005.

[7] J. Cho and H. Garcia-Molina. Synchronizing a database to improve freshness. In *SIGMOD*, 2000.

[8] G. Cormode and M. Garofalakis. Sketching streams through the net: distributed approximate query tracking. In *VLDB*, 2005.

[9] C. D. Cranor, T. Johnson, O. Spatscheck, and V. Shkapenyuk. Gigascope: A stream database for network applications. In *SIGMOD*, 2003.

[10] A. Deligiannakis, Y. Kotidis, and N. Roussopoulos. Hierarchical in-network data aggregation with quality guarantees. In *EDBT*, 2004.

[11] C. Estan and G. Varghese. New directions in traffic measurement and accounting. In *SIGCOMM*, 2002.

[12] M. J. Freedman and D. Mazires. Sloppy Hashing and Self-Organizing Clusters. In *IPTPS*, 2003.

[13] FreePastry. `http://freepastry.rice.edu`.

[14] P. K. Gummadi, R. Gummadi, S. D. Gribble, S. Ratnasamy, S. Shenker, and I. Stoica. The impact of dht routing geometry on resilience and proximity. In *SIGCOMM*, 2003.

[15] R. Gupta and K. Ramamritham. Optimized query planning of continuous aggregation queries in dynamic data dissemination networks. In *WWW*, pages 321–330, 2007.

[16] N. J. A. Harvey, M. B. Jones, S. Saroiu, M. Theimer, and A. Wolman. SkipNet: A Scalable Overlay Network with Practical Locality Properties. In *USITS*, March 2003.

[17] Y. Hu, D. M. Chiu, and J. C. Lui. Adaptive flow aggregation - a new solution for robust flow monitoring under security attacks. In *NOMS*, 2006.

[18] R. Huebsch, J. M. Hellerstein, N. Lanham, B. T. Loo, S. Shenker, and I. Stoica. Querying the Internet with PIER. In *VLDB*, 2003.

[19] N. Jain, L. Amini, H. Andrade, R. King, Y. Park, P. Selo, and C. Venkatramani. Design, implementation, and evaluation of the linear road benchmark on the stream processing core. In *SIGMOD*, 2006.

[20] N. Jain, D. Kit, P. Mahajan, P. Yalagandula, M. Dahlin, and Y. Zhang. PRISM: Precision-Integrated Scalable Monitoring (extended). Technical Report TR-07-15, UT Austin Department of Computer Sciences, 2007.

[21] N. Jain, D. Kit, P. Mahajan, P. Yalagandula, M. Dahlin, and Y. Zhang. Star: Self tuning aggregation for scalable monitoring. In *VLDB*, 2007.

[22] R. Keralapura, G. Cormode, and J. Ramamirtham. Communication-efficient distributed monitoring of thresholded counts. In *SIGMOD*, 2006.

[23] J. Kleinberg. Bursty and hierarchical structure in streams. In *KDD*, 2002.

[24] S. R. Madden, M. J. Franklin, J. M. Hellerstein, and W. Hong. TAG: a Tiny AGgregation Service for Ad-Hoc Sensor Networks. In *OSDI*, 2002.

[25] R. Mahajan, S. Bellovin, S. Floyd, J. Vern, and P. Scott. Controlling high bandwidth aggregates in the network. In *ACM SIGCOMM CCR*, 2001.

[26] C. Olston, J. Jiang, and J. Widom. Adaptive filters for continuous queries over distributed data streams. In *SIGMOD*, 2003.

[27] C. Olston and J. Widom. Offering a precision-performance tradeoff for aggregation queries over replicated data. In *VLDB*, Sept. 2000.

[28] C. Olston and J. Widom. Best-effort cache synchronization with source cooperation. In *SIGMOD*, 2002.

[29] C. G. Plaxton, R. Rajaraman, and A. W. Richa. Accessing Nearby Copies of Replicated Objects in a Distributed Environment. In *ACM SPAA*, 1997.

[30] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content Addressable Network. In *SIGCOMM*, 2001.

[31] A. Rowstron and P. Druschel. Pastry: Scalable, Distributed Object Location and Routing for Large-scale Peer-to-peer Systems. In *Middleware*, 2001.

[32] M. A. Shah, J. M. Hellerstein, and E. Brewer. Highly-available, fault-tolerant, parallel dataflows. In *SIGMOD*, 2004.

[33] A. S. Silberstein, R. Braynard, C. Ellis, K. Munagala, and J. Yang. A sampling-based approach to optimizing top-k queries in sensor networks. In *ICDE*, 2006.

[34] A. Singla, U. Ramachandran, and J. Hodgins. Temporal notions of synchronization and consistency in Beehive. In *Proc. SPAA*, 1997.

[35] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan. Chord: A scalable Peer-To-Peer lookup service for internet applications. In *ACM SIGCOMM*, 2001.

[36] N. Tatbul, U. Çetintemel, and S. Zdonik. Staying FIT: Efficient Load Shedding Techniques for Distributed Stream Processing. In *VLDB*, 2007.

[37] R. van Renesse, K. Birman, and W. Vogels. Astrolabe: A robust and scalable technology for distributed system monitoring, management, and data mining. *TOCS*, 2003.

[38] D. Veitch, S. Babu, and A. Pasztor. Robust synchronization of software clocks across the internet. In *IMC*, 2004.

[39] P. Yalagandula and M. Dahlin. A scalable distributed information management system. In *Proc SIGCOMM*, 2004.

[40] H. Yu and A. Vahdat. Design and evaluation of a conit-based continuous consistency model for replicated services. *TOCS*, 2002.