

STAR: Self-Tuning Aggregation for Scalable Monitoring

Navendu Jain, Dmitry Kit, Prince Mahajan, Praveen Yalagandula[†], Mike Dahlin, and Yin Zhang
Department of Computer Sciences [†]*Hewlett-Packard Labs*
University of Texas at Austin *Palo Alto, CA*

ABSTRACT

We present STAR, a self-tuning protocol that adaptively sets numeric precision constraints to accurately and efficiently answer continuous aggregate queries over distributed data streams. Adaptivity and approximation are essential for both robustness to varying workload characteristics and for scalability to large systems. In contrast to previous studies, we treat the problem as an optimization problem whose goal is to minimize the total communication load for a multi-level aggregation tree under a fixed error budget. Our hierarchical self-tuning algorithm, STAR, computes optimal error distribution and performs cost-benefit throttling to direct error slack to where it benefits the most. A key novel feature of STAR is that it takes into account the update rate and variance in the input data distribution in a principled manner. Our prototype implementation of STAR in a large-scale monitoring system provides (1) a new *distribution mechanism* that enables self-tuning error distribution and (2) an optimization to reduce communication overhead in a practical setting by carefully distributing the initial, default error budgets. Through extensive experiments performed on synthetic data simulations and a real network monitoring implementation, we show that STAR achieves significant performance benefits compared to existing approaches while still providing high accuracy and low resource utilization for dynamic workloads.

1. INTRODUCTION

This paper describes STAR, a *self-tuning* algorithm for adaptive setting of numeric precision constraints to process continuous aggregate queries in a large-scale monitoring system.

Scalable system monitoring is a fundamental abstraction for large-scale networked systems. It serves as a basic building block for applications such as network monitoring and management [8, 18, 41], financial applications [3], resource location [20, 40], efficient multicast [37], sensor networks [20, 40], resource management [40], and bandwidth provisioning [13]. To provide a real-time view of global system state for these monitoring and control applications, the central challenge for a monitoring system is scalability to process queries over multiple, continuous, rapid, time-varying data streams that generate updates for thousands or millions of *dynamic* attributes (e.g., per-flow or per-object state) spanning tens of thousands of nodes.

Recent studies [26, 28, 34, 37, 42] suggest that real-world applications often can tolerate some inaccuracy as long as the maximum error is bounded and that small amounts of approximation error can provide substantial bandwidth reductions. However, setting of a static error budget a priori is difficult when workloads are not known in advance and inefficient when workload characteristics change unpredictably over time. Therefore, for many real-world systems that require processing of long-running or continuous queries over data streams, it is important to consider *adaptive* ap-

proaches for query processing [2, 5, 7, 28].

A fundamental observation behind this work is that an adaptive algorithm for setting error budgets should be based on three key design principles for large-scale query processing :

- **Principled Approach:** First, to provide a general and flexible framework for adaptive error distribution for different workloads, we require a solution that does not depend on any a priori knowledge about the input data distribution. Rather, a self-tuning algorithm should be based on first principles and use the *workload* itself to guide the process of dynamically adjusting the error budgets.
- **Cost-Benefit Throttling:** Second, rather than continuously redistributing error budgets in pursuit of a perfect distribution, our algorithm explicitly determines when the current distribution is close enough to optimal that sending messages to redistribute allocations will likely cost more than it will ultimately save given the measured variability of the workload. For example, in our network monitoring service for detecting elephant flows [13], we track bandwidth for tens of thousands of flows, but the vast majority of these flows are mice that produce so few updates that further redistribution of the error budgets is not useful. Avoiding fruitless optimization in such cases significantly improves scalability in systems with tens of thousands of attributes.
- **Aggregation Hierarchy:** Third, to provide a global view of the system requires processing and aggregating data from distributed data streams in real-time. In such an environment, a hierarchical decentralized query processing infrastructure provides an attractive solution to minimize communication/computation costs for both scalability and load balancing. Further, in a hierarchical aggregation tree, the internal nodes can not only split the error budget among their children but may also retain some local error budget to prevent updates received from children from being propagated further up the tree e.g., merging of updates that cancel out each other.

Unfortunately, existing approaches do not satisfy these requirements. On one hand, protocols such as adaptive filters [4, 28] and adaptive thresholded counts [24] can effectively reduce communication overhead given a fixed error budget for flat (2-tier) topologies, but they offer neither scalability to a large number of nodes and attributes nor functionality of in-network aggregation. On the other hand, existing hierarchical protocols either assign a static error budget [27] that cannot adapt to changing workloads, or periodically shrink error thresholds [11] at each node to create redistribution error budget that incurs high load for monitoring a large number of attributes. Finally, although the previous solutions have an intuitive appeal, they lack a rigorous problem formulation leading to solutions that are difficult to compare against the global optimization problem of minimizing the total communication overhead.

To address these challenges, STAR design focuses on pro-

viding three key properties for adaptively setting error thresholds in hierarchical aggregation to provide high performance, scalability, and adaptivity.

- **High Performance:** To compute optimal error assignments, STAR uses an informed mathematical model to formulate an optimization problem whose goal is to minimize the global communication load in an aggregation tree given a fixed total budget. This model provides a closed-form, optimal solution for adaptive setting of error budgets using only local and aggregated information at each node in the tree. Given the optimal error budgets, STAR performs cost-benefit throttling to balance the tradeoff between the cost for redistributing error budgets and the expected benefits. Our experimental results show that self-tuning distribution of error budgets can reduce monitoring costs by up to a factor of five over previous approaches.
- **Scalability:** STAR builds on recent work that uses distributed hash tables (DHTs) to construct scalable, load-balanced forests of self-organizing aggregation trees [40, 14, 6, 30]. Scalability to tens of thousands of nodes and millions of attributes is achieved by mapping different attributes to different trees. In this framework of forest of aggregation trees, STAR’s self-tuning algorithm directs error slack to where it is most needed.
- **Convergence:** STAR computes optimal error budgets and performs cost-benefit analysis to continually adjust to dynamic workloads. But since STAR’s self-tuning algorithm adapts its solution as input workload changes, it is difficult to qualify its convergence properties. Nonetheless, under stable input data distributions, we show that STAR converges to a solution with desirable properties. Further, STAR balances the speed of adaptivity and robustness to workload fluctuations.

We study the performance of our algorithm through both synthetic data simulations and a prototype implementation on our SDIMS aggregation system [40] based on FreePastry [15]. Experience with a Distributed Heavy Hitter detection (DHH) application built on STAR illustrate how explicitly managing imprecision can qualitatively enhance a monitoring service. Our experimental results show the improved performance and scalability benefits: for the DHH application, small amounts of imprecision drastically reduce monitoring load. For example, a 10% error budget allows us to reduce network load by an order of magnitude compared to the uniform allocation policy. Further, for 90:10 skewness in attributes (i.e., 10% heavy hitters) in terms of input load, self-tuning gives more than an order of magnitude better performance than both uniform error allocation and approaches that do not perform cost-benefit throttling.

This paper makes three key contributions. First, we present STAR, the first self-tuning algorithm for scalable aggregation that computes optimal error distribution and performs cost-benefit throttling for large-scale system monitoring. Second, we provide a scalable implementation of STAR in our SDIMS monitoring system. Our implementation provides a new *distribution abstraction*, a dual mechanism to the traditional bottom-up tree based aggregation, that enables *self-tuning* error distribution top-down along the aggregation tree to reduce communication load. Further, it provides an important optimization that reduces communication overhead in a practical setting by carefully distributing the initial, default error budgets. Third, our evalua-

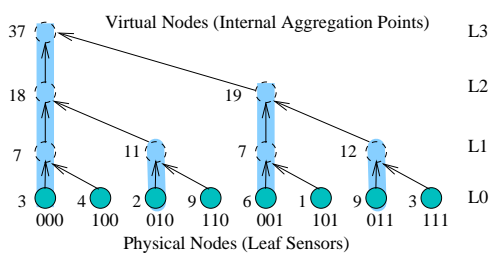


Figure 1: The aggregation tree for key 000 in an eight node system. Also shown are the aggregate values for a simple SUM() aggregation function.

tion demonstrates that adaptive error distribution is vital for enabling scalable aggregation: a system that performs self-tuning of error budgets can significantly reduce communication overheads.

The rest of this paper is organized as follows. Section 2 provides background description of SDIMS [40], a scalable DHT-based aggregation system that underlies STAR. Section 3 describes the mechanism and the STAR self-tuning algorithm for adaptively setting arithmetic imprecision for reducing monitoring overhead. Section 4 presents the implementation of STAR in our SDIMS aggregation system, policies for maintaining precision of query results under failures, and policies for initializing the error budgets. Section 5 presents the experimental evaluation of our system. Finally, Section 6 discusses related work, and Section 7 provides conclusions.

2. BACKGROUND

SDIMS builds on two ongoing research efforts for scalable monitoring: aggregation and DHT-based aggregation.

2.1 Aggregation.

Aggregation is a fundamental abstraction for scalable monitoring [6, 14, 20, 30, 37, 40] because it allows applications to access summary views of global information and detailed views of rare events and nearby information.

SDIMS’s aggregation abstraction defines a tree spanning all nodes in the system. As Figure 1 illustrates, each physical node in the system is a leaf and each subtree represents a logical group of nodes. Note that 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 on a LAN in the CS department) [17, 40]. An internal non-leaf node, which we call a *virtual node*, is simulated by one or more physical nodes at the leaves of the subtree rooted at the virtual node.

The tree-based aggregation in the SDIMS framework is defined in terms of an aggregation function installed at all the nodes in the tree. Each leaf node (physical sensor) inserts or modifies its local value for an *attribute* defined as an {attribute type, attribute name} pair which is recursively aggregated up the tree. For each level- i subtree T_i in an aggregation tree, SDIMS defines an *aggregate value* $V_{i,attr}$ for each attribute: for a (physical) leaf node T_0 at level 0, $V_{0,attr}$ is the locally stored value for the attribute or NULL if no matching tuple exists. The aggregate value for a level- i subtree T_i is the result returned by the aggregation function

computed across the aggregate values of T_i 's children. Figure 1, for example, illustrates the computation of a simple SUM aggregate.

2.2 DHT-based aggregation.

SDIMS leverages DHTs [30, 31, 32, 36, 44] to construct a forest of aggregation trees and maps different attributes to different trees [6, 14, 30, 33, 40] for scalability and load balancing. DHT systems assign a long (e.g., 160 bits), random ID to each node and define a routing algorithm to send a request for key k to a node $root_k$ such that the union of paths from all nodes forms a tree $DHTtree_k$ rooted at the node $root_k$. By aggregating an attribute with key $k = \text{hash}(\text{attribute})$ along the aggregation tree corresponding to $DHTtree_k$, different attributes are load balanced across different trees. Studies suggest that this approach can provide aggregation that scales to large numbers of nodes and attributes [40, 33, 30, 14, 6].

2.3 Example Application

Aggregation is a building block for many distributed applications such as network management [41], service placement [16], sensor monitoring and control [26], multicast tree construction [37], and naming and request routing [9]. In this paper, we focus on a case-study example: a distributed heavy hitter detection service.

2.3.1 Heavy Hitter detection

Our case-study application is identifying heavy hitters¹ in a distributed system—for example, the top 10 IPs that account for a significant fraction of total incoming traffic in the last 10 minutes [13]. The key challenge for this distributed query is scalability for aggregating per-flow statistics for tens of thousands to millions of concurrent flows in real-time. For example, a subset of the Abilene [1] traces used in our experiments include up to 80 thousand flows that send about 25 million updates per hour.

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 IP from all nodes in the system using SUM as the aggregation function and $\text{hash}(\text{HH-Step1}, \text{destIP})$ as the key. For example, tuple (H = $\text{hash}(\text{HH-Step1}, 128.82.121.7)$, 700 KB) at the root of the aggregation tree T_H indicates that a total of 700 KB of data was received for 128.82.121.7 across all vantage points during the last time window. In the second step, we feed these aggregated total bandwidths for each destination IP into a SELECT-TOP-10 aggregation with key $\text{hash}(\text{HH-Step2}, \text{TOP-10})$ to identify the TOP-10 heavy hitters among all flows.

Although there exist other centralized monitoring services, in Section 5 we show that using our STAR self-tuning algorithm in the SDIMS aggregation system, we can monitor a large number of attributes at much finer time scales while incurring significantly lower network costs.

¹Note that the standard definition of a heavy hitter is an entity that accounts for at least a specified proportion of the total activity measured in terms of number of packets, bytes, connections etc [13]. We use a slightly different definition of “heavy hitters” to denote flows whose bandwidth is greater than a specified fraction threshold of the maximum flow value.

3. STAR DESIGN

Arithmetic imprecision (AI) bounds the numeric difference between a reported value of an attribute and its true value [29, 43]. For example, a 10% AI bound ensures that the reported value either underestimates or overestimates the true value by at most 10%.

When applications do not need exact answers and data values do not fluctuate wildly, arithmetic imprecision can greatly reduce the monitoring load by allowing caching to filter small changes in aggregated values. Furthermore, for applications like distributed heavy hitter monitoring, arithmetic imprecision can completely filter out updates for most “mice” flows.

We first describe the basic mechanism for enforcing AI for each aggregation subtree in the system. Then we describe how our system uses a self-tuning algorithm to address the policy question of distributing an AI budget across subtrees to minimize system load.

3.1 Mechanism

To enforce AI, each aggregation subtree T for an attribute has an error budget δ_T which 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 δ_{self} and its children δ_c (with $\delta_T \geq \delta_{self} + \sum_{c \in \text{children}} \delta_c$), and the children recursively do the same. Here we present the AI mechanism for the SUM aggregate since they are likely to be more common in practice in network monitoring and financial applications; other standard aggregation functions (e.g., MAX, MIN, AVG) are similar [23].

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 δ_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

$$\begin{aligned} V_{min} &\leq \sum_{c \in \text{children}} V_{min}^c \\ V_{max} &\geq \sum_{c \in \text{children}} V_{max}^c \end{aligned} \quad (1)$$

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

Note the trade-off in splitting δ_T between δ_{self} and δ_c . Large δ_c allows children to filter updates before they reach a node. Conversely, by setting $\delta_{self} > 0$, a node can set $V_{min} < \sum V_{min}^c$, set $V_{max} > \sum V_{max}^c$, or both to avoid further propagating some updates it receives from its children.

SDIMS maintains per-attribute δ values so that different attributes with different error requirements and different update patterns can use different δ budgets in different subtrees. SDIMS implements this mechanism by defining a *distribution function*; just as an attribute type’s aggregation function specifies how aggregate values are aggregated from children, an attribute type’s distribution value specifies how δ budgets are distributed among children and δ_{self} .

3.2 Policy Decisions

Given these mechanisms, there is a plenty of flexibility to (i) set δ_{root} to an appropriate value for each attribute, and (ii) compute V_{min} and V_{max} when updating a parent.

Setting δ_{root} : Note that the aggregation queries can set the root error budget δ_{root} to any non-negative value. For some applications, an absolute constant value may be known a priori (e.g., count the number of connections per second ± 10 at port 1433.) For other applications, it may be appropriate to set the tolerance based on measured behavior of the aggregate in question (e.g., set δ_{root} for an attribute to be at most 10% of the maximum value observed) or the measurements of a set of aggregates (e.g., in our heavy hitter application, set δ_{root} for each flow to be at most 1% of the bandwidth of the largest flow measured in the system). Our algorithm supports all of these approaches by allowing new absolute δ_{root} values to be introduced at any time and then distributed down the tree via a distribution function. We have prototyped systems that use each of these three policies.

Computing $[V_{min}, V_{max}]$: When either $\sum_c V_{min}^c$ or $\sum_c V_{max}^c$ goes outside of the last $[V_{min}, V_{max}]$ that was reported to the parent, a node needs to report a new range. Given a δ_{self} budget at an internal node, we have some flexibility on how to center the $[V_{min}, V_{max}]$ range. Our approach is to adopt a per-aggregation-function range policy that reports $V_{min} = (\sum_c V_{min}^c) - bias * \delta_{self}$ and $V_{max} = (\sum_c V_{max}^c) + (1 - bias) * \delta_{self}$ to the parent. The *bias* parameter can be set as follows:

- *bias* ≈ 0.5 if inputs expected to be roughly stationary
- *bias* ≈ 0 if inputs expected to be generally increasing
- *bias* ≈ 1 if inputs expected to be generally decreasing

For example, suppose a node with total δ_T of 10 and δ_{self} of 3 has two children reporting ($[V_{min}^c, V_{max}^c]$) of [1, 2] and [2, 8], respectively, and reports [0, 10] to its parent. Then, the first child reports a new range [10, 11], so the node must report to its parent a range that includes [12, 19]. If *bias* = 0.5, then report to parent [10.5, 20.5] to filter out small deviation around the current position. Conversely, if *bias* = 0, report [12, 22] to filter out the maximal number of updates of increasing values.

3.3 Self-tuning Error Budgets

The key AI policy question is how to divide a given error budget δ_{root} across the nodes in an aggregation tree.

A simple approach is a static policy that divides the error budget uniformly among all the children. E.g., a node with budget δ_T could set $\delta_{self} = 0.1\delta_T$ and then divide the remaining $0.9\delta_T$ evenly among its children. Although this approach is simple, it is likely to be inefficient because different aggregation subtrees may experience different loads.

To make cost/accuracy tradeoffs *self-tuning*, we provide an adaptive algorithm. The high-level idea is simple: increase δ for nodes with high load and standard deviation but low δ (relative to other nodes); decrease δ for nodes with low load and low standard deviation but high δ . Next, we address the problem of optimal distribution of error budgets for a 2-tier (one-level) tree and later extend it as a general approach for a hierarchical aggregation tree.

3.3.1 One-level tree

Quantify AI Filtering Gain: In order to derive the optimal distribution of error budgets among different nodes, we need a simple way of quantifying the amount of load reduction that can be achieved when a given error budget is used for AI filtering.

Intuitively, the AI 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 2, if the allocated error budget δ_i at node i is much smaller than the standard deviation σ_i of the underlying data distribution, δ_i is unlikely to filter many data updates. Meanwhile, if δ_i is above σ_i , we would expect the load to decrease quickly as δ_i increases.

In order to quantify the tradeoff between load and error budget, one possibility is to compute the entire tradeoff curve as shown in Figure 2. However, doing so imposes several difficulties. First, it is in general difficult to compute the tradeoff curve without a priori knowledge about the underlying data distribution. Second, maintaining the entire tradeoff curve becomes expensive when there are a large number of attributes and nodes. Finally, it is not easy to optimize the distribution of error budgets among different nodes based on the tradeoff curves.

To overcome these difficulties, we develop a simple metric in STAR to capture the tradeoff between load and error budget. Our metric is motivated by the Chebyshev's inequality in probability theory, 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 μ and variance σ^2 . The Chebyshev's inequality states that for any $k \geq 0$, the probability of the event

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

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

$$Pr(|X - \mu| \geq k\sigma_i) \leq \frac{1}{k^2} = \frac{\sigma_i^2}{\delta_i^2} \quad (3)$$

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 δ_i is unlikely to filter any data updates (Figure 2). In this case, Equation 3 provides only a weak bound on message cost: the probability that each incoming update will trigger an outgoing message is upper bounded by 1. However, if $\delta_i \geq k\sigma_i$ for any $k \geq 1$, the fraction of unfiltered updates is bounded by $\frac{\sigma_i^2}{\delta_i^2}$.

In general, given the input update rate u_i for child i with error budget δ_i , the expected cost for node i per unit time is:

$$\text{MIN} \left(1, \frac{\sigma_i^2}{\delta_i^2} \right) * u_i \quad (4)$$

Compute Optimal Error Budget: To derive the optimal error distribution at each node, we can formulate an optimization problem of minimizing the total incoming network load at root R under a fixed total AI budget δ_T i.e.,

$$\begin{aligned} \text{MIN} \quad & \sum_{i \in \text{child}(R)} \frac{\sigma_i^2 * u_i}{(\delta_i^{\text{opt}})^2} \\ \text{s.t.} \quad & \sum_{i \in \text{child}(R)} \delta_i^{\text{opt}} = \delta_T \end{aligned} \quad (5)$$

Using Lagrange multipliers yields a closed-form and computationally inexpensive optimal solution [23]:

$$\delta_i^{\text{opt}} = \delta_T * \frac{\sqrt[3]{\sigma_i^2 * u_i}}{\sum_{i \in \text{child}(R)} \sqrt[3]{\sigma_i^2 * u_i}} \quad (6)$$

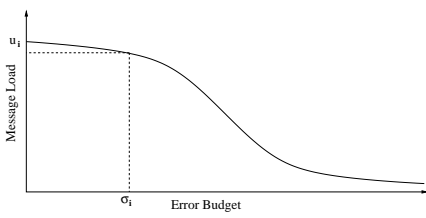


Figure 2: Expected load vs error budget of a node.

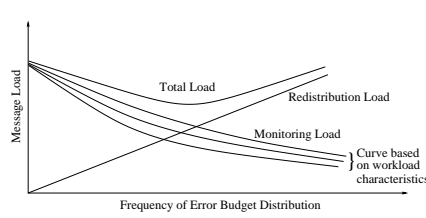


Figure 3: Cost-benefit Analysis.

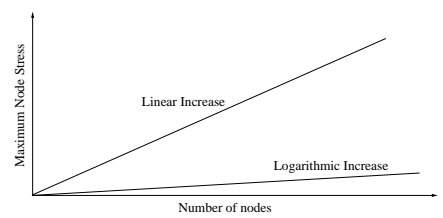


Figure 4: Maximum Node stress: 1-level vs hierarchical tree

The above optimal error assignment assumes that for all nodes, the expected cost per update is equal to $\frac{\sigma_i^2}{\delta_i}$ based on Equation 3 i.e., $\sigma_i \leq \delta_i$. However, for nodes with high σ_i relative to error budget δ_i , it is highly likely that an update will be sent to the root for each incoming message. We term these nodes as *volatile* nodes [11] since they may not yield a significant benefit in spite of being allocated a large fraction of the error budget (Equation 6).

To account for volatile nodes, we apply an iterative algorithm that determines the largest volatile node j at each step, recomputes Equation 6 for all remaining children assuming j is absent. A node j is labeled volatile if $\frac{\sigma_j}{\delta_j^{opt}} \geq 1$ i.e., the standard deviation is larger than the optimal error budget (under fixed total budget) corresponding to Equation 3 and the ratio $\frac{\sigma_j}{\delta_j^{opt}}$ is maximal among all remaining children. If no such j exists, the procedure terminates giving the optimal AI budgets for each node; all volatile nodes get zero budget since any non-zero budget will not effectively filter their updates. Note that for our DHT-based aggregation trees, the fan-in for a node is typically 16 (i.e., 4-bit correction per hop) so the iterative algorithm runs in constant time (at most 16 times).

Relaxation: A self-tuning algorithm that adapts too rapidly may react inappropriately to transient situations. Therefore, we next apply exponential smoothing to compute the new error budget δ_i^{new} for each node as the weighted average of new error budget (δ_i^{opt}) and previous budget (δ_i).

$$\delta_i^{new} = \alpha \delta_i^{opt} + (1 - \alpha) \delta_i \quad (7)$$

where $\alpha = 0.05$.

Cost-Benefit Throttling: Finally, root R needs to send messages to its children to rebalance the error budget. Therefore, there is a tradeoff between the error budget redistribution overhead and the AI filtering gain (as illustrated in Figure 3). A naive rebalancing algorithm that ignores such tradeoff could easily spend more network messages redistributing δ s than it saves by filtering updates. Limiting redistribution overhead is a particular concern for applications like distributed heavy hitter 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, the root node computes a *charge* metric for each child subtree c , which estimates the number of extra messages sent by c due to sub-optimal δ :

$$Charge_c = (T_{curr} - T_{adjust}) * (M_c - M_c^{new})$$

where $M_c = \frac{\sigma_i^2 * u_i}{\delta_i^2}$, $M_c^{new} = \frac{\sigma_i^2 * u_i}{(\delta_i^{new})^2}$, and T_{adjust} is the last time δ was adjusted at R for child c . Notice that a subtree's

charge will be large if (a) there is a large load imbalance (e.g., $M_c - M_c^{new}$ is large) or (b) there is a stable, long-lasting imbalance (e.g., $T_{curr} - T_{adjust}$ is large.)

We only send messages to redistribute deltas when doing so is likely to save at least k messages (i.e., if $charge_c > k$). To ensure the invariant that $\delta_T \geq \delta_{self} + \sum_c \delta_c$, we make this adjustment in two steps. First, we replenish δ_{self} from the child whose δ_c is the farthest above δ_c^{new} by ordering c to reduce δ_c by $\text{Min}(0.1 \delta_c, \delta_c - \delta_c^{new})$. Second, we loan some of the δ_{self} budget to the node c that has accumulated the largest charge by incrementing c 's budget by $\text{Min}(0.1 \delta_c, \delta_c^{new} - \delta_c, \max(0.1 \delta_{self}, \delta_{self} - \delta_{self}^{new}))$.

A node responds to a request from its parent to update δ_T using a similar approach.

Note that in a practical setting, a parent node should first reclaim error budget from children and then give budget to other children nodes as guided by the self-tuning algorithm.

3.3.2 Multi-level trees

For large-scale multi-level trees, we extend our basic algorithm for a one-level tree to a distributed algorithm for a multi-level aggregation hierarchy to reduce maximum node stress (Figure 4) and reduce communication load due to internal nodes that not only split δ_c among their children c but may also retain δ_{self} to help prevent updates received from children from being propagated further up the tree. Second, in order to scale to the tens of thousands of attributes required by per-flow network monitoring services, our algorithm performs cost/benefit throttling. In particular, we use a competitive analysis to ensure that rebalancing δ accounts for only a fraction of the system's total message load by only redistributing δ to correct large imbalances or to correct stable, long-duration (though perhaps small) imbalances.

At any internal node in the aggregation tree, the self-tuning algorithm works as follows:

1. Estimate optimal distribution of δ_T across δ_{self} and δ_c . Each 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 (σ_p) of updates received from its children. Note that u_c, σ_c reports are accumulated by child c until they can be piggy-backed on an update message to its parent.

Given this information, each parent node n computes the optimal values δ_v^{opt} for each child v 's underlying subtree that minimizes the total system load in the entire subtree rooted at n . We apply Equation 6 by viewing each child v as representing two individual nodes: (1) v itself (as a data source) with update rate u_v and standard deviation σ_v and (2) a node representing the subtree rooted at v . Figure 5 illustrates this local self-tuning view: when any internal node

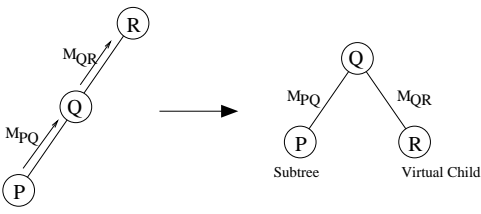


Figure 5: Local self-tuning view considers both incoming and outgoing bandwidth.

computes optimal budgets, it aims to minimize both incoming messages received as well as outgoing messages to parent (by showing parent-link as a virtual child) i.e., minimizing global communication load from a local perspective.

Given this view model, we define $loadFactor$ for a node v as $\sqrt[3]{\sigma_v^2 * u_v}$. Recursively, we can define $AccLoadFactor$ for a subtree rooted at node v as:

$$AccLoadFactor_v = \begin{cases} LoadFactor_v & (v \text{ is a leaf node}) \\ LoadFactor_v + \sum_{j \in child(v)} AccLoadFactor_j & \\ (otherwise) & \end{cases}$$

Next, we estimate the optimal error budget for v 's subtree ($v \in child(n)$) as follows:

$$\delta_v^{opt} = \delta_T * \frac{AccLoadFactor_v}{AccLoadFactor_n} \quad (8)$$

Equation 8 is globally optimal since it virtually maps a multi-level hierarchy into a one-level tree (as illustrated in Figure 6) and in this transformed view, computes the optimal error budget for each node.

To account for volatile nodes, we apply a similar iterative algorithm as in the one-level tree case to determine the largest volatile node i ($i \in child(n)$) at each step, recompute Equation 8 for all remaining children, and so on. The condition to check whether node i is volatile remains same as: $\frac{\sigma_i}{\delta_{i(leaf)}^{opt}} \geq 1$ and $\frac{\sigma_i}{\delta_{i(leaf)}^{opt}} \geq \frac{\sigma_j}{\delta_{j(leaf)}^{opt}} \forall j$ where $\delta_{i(leaf)}^{opt}$ is computed as:

$$\delta_{i(leaf)}^{opt} = \delta_T * \frac{LoadFactor_i}{AccLoadFactor_v} \quad (9)$$

2. *Relaxation: adaptive adjustment of deltas.*

$$\delta_i^{new} = \alpha \delta_i^{opt} + (1 - \alpha) \delta_i$$

where $\alpha = 0.05$.

3. *Redistribute deltas iff the expected benefit exceeds the redistribution overhead.*

To do cost-benefit throttling, We recursively apply the formula for one-level tree to compute the cumulative charge for a subtree rooted at node v : the following difference: for computing $M_i - M_i^{new}$,

$$AccCharge_v = \begin{cases} Charge_v & (v \text{ is a leaf node}) \\ Charge_v + \sum_{j \in child(v)} AccCharge_j & \\ (otherwise) & \end{cases}$$

$$M_i = \left(\frac{\sigma_i^2 * u_i}{\delta_i^2} + \sum_{j \in child(i)} \frac{\sigma_j^2 * u_j}{\delta_j^2} \right)$$

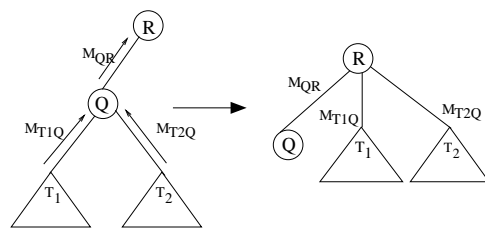


Figure 6: Global self-tuning view collapses a multi-level hierarchy to a one-level tree.

and

$$M_i^{new} = \left(\frac{\sigma_i^2 * u_i}{(\delta_i^{new})^2} + \sum_{j \in child(i)} \frac{\sigma_j^2 * u_j}{(\delta_j^{new})^2} \right)$$

Note that the terms $AccLoadFactor_v$ for computing Equation 8 and $AccCharge_v$ can be computed using a SUM aggregation function, and are piggybacked on updates sent by children to their parents in our implementation.

While beyond the scope of this paper, our algorithms can be extended in a straightforward manner to adaptively assign error budgets for multiple queries involving overlapping sets of data objects.

4. IMPLEMENTATION

In this section, we describe three important design aspects of implementing STAR in our SDIMS prototype. First, we present a new *distribution abstraction* to provide the functionality of distributing AI error budgets in an aggregation tree. Second, we discuss different policies to maintain AI error precision in query results under failures. Finally, we describe two practical optimizations in our prototype implementation to reduce communication load for large-scale system monitoring.

4.1 Distribution Abstraction

A distribution abstraction is a dual mechanism to the traditional bottom-up tree based aggregation that enables an update to be *distributed* top-down along the aggregation tree.

The basic mechanism of a distribution function at a node is to receive its list of children for an attribute, a (possibly null) Distribution message received from its parent, and return a set of messages destined for a subset of its children and itself. For self-tuning AI, a distribution function implements the STAR algorithm that takes an AI error budget from the parent and distributes it down to its children.

In the SDIMS framework, each node implements an Aggregation Management Layer (AML) that maintains attribute tuples, performs aggregations, stores and propagates aggregate values [40]. On receiving a Distribution message, the AML layer invokes the distribution function for the given attribute and level in the aggregation tree. The logic of how to process these messages is implemented by the application. For example, for self-tuning AI, the STAR protocol generates distribution messages to either allocate more error budget to a child subtree or decrease the error budget of a child subtree.

The AML layer provides communication between an aggregation function and a distribution function by passing local messages. Conceptually, a distribution function ex-

tends an aggregation function by allowing the flexibility of defining multiple distribution functions (policies) for a given aggregation function.

In SDIMS, a store of (attribute type, attribute name, value) tuples is termed as Management Information Base (MIB). For hierarchical aggregation of a given attribute, each node stores *child MIBs* received from children and a *reduction MIB* containing locally aggregated values across the child MIBs. To provide the distribution abstraction, we implemented a *distribution MIB* that stores as value an object used for maintaining application-defined state for adaptive error budget distribution. In STAR, a distribution MIB also holds a local copy of load information received from child MIBs.

Given this abstraction, it is simple to implement STAR’s self-tuning adaptation of error budgets in an aggregation tree. For example, in STAR

- (a) A parent can increase δ_{self} by sending a distribution message to a child c to reduce its subtree budget, δ_c
- (b) A parent can increase δ_c for a subtree child c by reducing its own δ_{self} , and
- (c) Filter small changes: a parent only changes a δ assignment if doing so is likely to save more messages than rebalancing the δ s costs.

The policy decision of when the self-tuning function gets called can be based on either a periodic timer, processing a threshold number of messages, or simply on-demand.

4.2 Robustness

To handle node failures and disconnections, our STAR implementation currently provides a simple policy to maintain the AI error precision in query results under churn. On each invocation of the distribution function, it receives the current child MIB set from the AML layer. If the new child MIB set is inconsistent with the previous child MIB snapshot maintained by the distribution function, it takes corrective action as follows:

- **On a new child event:** inserts a new entry in the distribution MIB for the new-child, assigns an error budget to that child subtree, and sends it to that child for distribution in its subtree.
- **On a dead child event:** garbage collects the child state and reclaims all AI error budget previously assigned to that child subtree.
- **On a new parent event:** Resets the AI error budget to zero as the new parent will allocate a new error budget.

Assuming that the error bounds were in a state where all precision constraints are satisfied prior to a failure, the temporary lost error due to the failure of a child only improves precision, thus no precision constraints can become violated. For a new child, it receives a fraction of the error budget allocation so the correctness still holds. Finally, on a new parent event, the error budget at its children will be reset.

Though this policy is simple and conservative, it always provides correctness that the reported AI error precision in query results is satisfied. Under this policy, however, a single failure might incur high communication overhead e.g., if a whole subtree moves to a new parent, the entire AI budget is lost resulting in every leaf update being propagated up in the tree until the new parent sends distribution messages to reassign the AI error budgets in the affected subtree.

Alternatively, policies based on exploiting consistency versus performance tradeoffs can be used. One such policy is when a child gets connected to a new parent, it can keep reporting aggregate values to the new parent with the precision error assigned by the last parent; the self-tuning algorithm then continually adjusts the previous subtree AI budget to converge with the new allocation over time.

Another policy would be to allow the invariant that a subtree meets the AI bound to violate under periods of re-configuration. Instead, a new parent can send down a target budget and aggregate up the actual AI error bound reported by child’s subtree. This policy is simple for both this failure case and the common case of adaptively redistributing error budgets. We leave the empirical comparison of different policies that trade strict consistency for performance as future work.

4.3 Optimizing for Scalability

With the STAR algorithm described in Section 3, the error budgets can be adaptively set in an aggregation hierarchy to minimize the communication overhead. However, since we are interested in large-scale system monitoring, in this section we present an optimization for setting initial error budgets that complement adaptive error settings to further reduce the network load in a real-world setting.

In our DHH application, we want to maintain accurate information for only the heavy hitter flows that generate a significant fraction of the total traffic. For mice flows that seldom send updates, we allow them to get culled at the leaves and maintain conservative information on their aggregate value.

Since SDIMS is an event-driven system, for error distribution, an initial update needs to be propagated to the root of an aggregation tree to distribute the error budget among all the nodes in that tree. Therefore, the initial cost of redistribution budget would be $O(\log N)$ for a mice update to reach the root and $O(N)$ for error distribution. However, for mice flows that only send a few updates, this cost will be significantly higher than benefits of filtering mice. Therefore for scalability, the key question is how to set the initial error budgets so that the mice flows never generate any updates.

One simple yet expensive policy is to keep the entire budget at the root and perform on-demand error distribution. Another policy is to give root share of zero and uniformly divide the entire budget at all the leaf nodes. In general, there is a continuum of policies that can be defined on initially dividing the total budget δ_T among the root and N leaf nodes. These set of policies can be expressed as:

$$AI_{Root} = Root_{share} * \delta_T$$

$$AI_{Leaf} = \frac{1 - AI_{Root}}{N}$$

For example, if $Root_{share} = 0.5$, then the root gets $\frac{\delta_T}{2}$ and each leaf gets $\frac{\delta_T}{2*N}$. To implement this policy, on receiving an update, a leaf node checks if it already has a local AI error budget. If not yet, it performs a lookup for AI_{Leaf} to filter the update. If the local budget is insufficient, the leaf sends the update to its parent which in turn applies AI error filtering. If that update reaches the root, the root initiates error distribution of its AI_{Root} error budget across its tree. Given this setting, we can cull a large fraction of the mice flows at the leaves preventing their updates from reaching the root. We show the performance benefit of using this

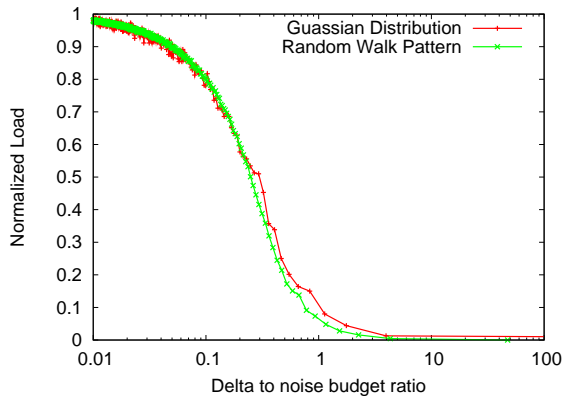


Figure 7: Normalized load vs. noise of synthetic workload for a fixed AI error budget. If noise $<$ AI, a majority of updates get filtered. The x-axis is on a log scale.

optimization in Section 5.

5. EXPERIMENTAL EVALUATION

We have implemented a prototype of STAR in our SDIMS monitoring framework [40] on top of FreePastry [32] and performed large-scale simulation experiments and a real-world monitoring system on two real networks: 120 node instances mapped on 30 physical machines in the department Condor cluster and the same setup on 30 machines in the Emulab [39] testbed.

Our experiments characterize the performance and scalability of the self-tuning AI for the distributed heavy hitters (DHH) application. First, we quantify the reduction in monitoring overheads due to self-tuning AI using simulations. Second, we evaluate the performance benefits of our optimization in a real world monitoring implementation. Finally, we investigate the reduction in communication load achieved by STAR for the DHH application using our prototype implementation. In summary, our experimental results show that STAR is an effective substrate for scalable monitoring: introducing small amounts of AI error and adaptivity using self-tuning AI significantly reduces monitoring load.

5.1 Simulation Experiments

To generalize the trade-off between AI error budget and monitoring load, we evaluate via simulations under what conditions is AI error budget and self-tuning the AI error budget effective. In all experiments, all active sensor are at the leaf nodes of the aggregation tree. Each sensor generates a data value every time unit (round) for two sets of synthetic workloads for 100,000 rounds: (1) a Gaussian distribution with standard deviation 1 and mean 0, and (2) a random walk pattern in which the value either increases or decreases by an amount sampled uniformly from $[0.5, 1.5]$.

We first investigate under what conditions is AI error budget effective. Figure 7 shows the simulation results for a 4-level degree-6 aggregation tree with 1296 leaf nodes for the two data distributions under uniform static error distribution. The x-axis denotes the ratio of total AI budget to total noise induced by the leaf sensors and the y-axis shows the total message load normalized with respect to zero AI error budget. We observe that when noise is small compared to

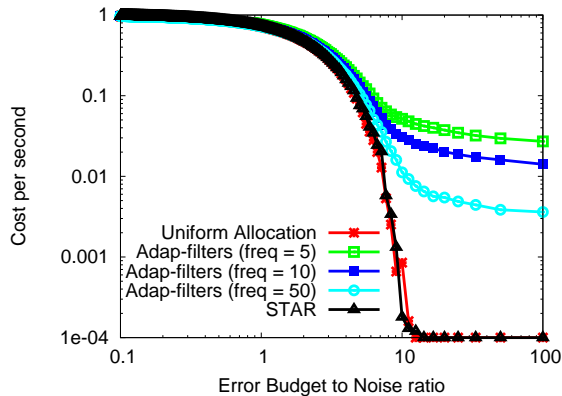


Figure 8: Performance benefits due to cost-benefit throttling. Load vs. precision error (log scale) for 10-node 1-level tree, random walk data.

the error budget, there is about an order of magnitude load reduction as the majority of updates are filtered. But, as expected, when noise is large compared to the error budget, the load asymptotically approaches the unfiltered load with $AI = 0$. The random walk pattern allows almost perfect culling of updates for small amounts of noise whereas for the Gaussian distribution, there is a small yet a finite probability for data values to deviate arbitrarily from their previously reported range.

Next, we quantify the cost of periodic bound shrinking in previous approaches [11, 28] compared with STAR that performs cost-benefit throttling. In our experiments, the following configuration gave the best results for Adaptive-filters [28]: shrink percentage = 5%, high self-tuning frequency, and distributing error budgets to a small set (e.g., 10-15%) of nodes with the highest burden scores, where burden is the ratio of load to error budget. These observations are consistent with previous work [11].

Figure 8 shows the performance results of uniform allocation, Adaptive-filters [28], and STAR for a 10 node 1-level tree for 1 attribute using a random walk pattern with same step size at each node. In this case, the uniform error allocation would be close to the optimal setting. We observe that when the error budget exceeds noise, Adaptive-filters incurs a constant error redistribution cost that is proportional to frequency of self-tuning. As discussed in Section 1, for per-flow monitoring services that require tracking millions of attributes each having small but non-zero noise, the corresponding communication overhead would be very expensive for Adaptive-filters. STAR, however, performs cost-benefit throttling and doesn't redistribute error when the corresponding gain is negligible.

We now characterize the higher performance benefits of STAR compared to other approaches as skewness in a workload increases. In Figure 9, the three figures show the communication load vs. error budget to noise ratio for the following skewness settings: (a) 20:80% (b) 50:50% (c) 90:10%. For example, the 20:80% skewness represents that only 20% of nodes have zero noise and the remaining 80% flows have a large noise. In this case, since only a small fraction of the nodes are stable, both STAR and Adaptive-filters can only reclaim 20% total error budget from the zero noise sources and distribute it to noisy sources to cull their up-

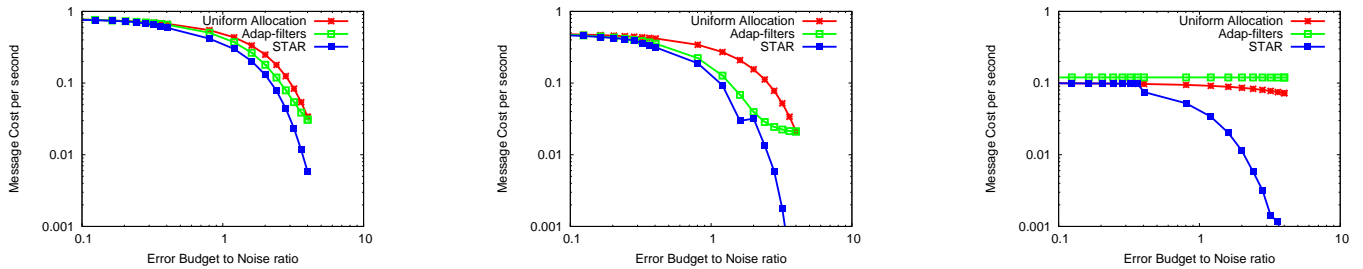


Figure 9: STAR provides higher performance benefits as skewness in a workload increases. The three figures show load vs. error budget to noise ratio for different skewness settings (a) 20:80% (b) 50:50% (c) 90:10%.

dates. STAR has a constant difference from Adaptive-filters since the latter does a periodic shrinking of error bound at each node. For the 50:50 case, both self-tuning algorithms can claim 50% of total budget compared to uniform allocation and give it to noisy sources. However, even when the optimal configuration (error budget large compared to noise) is reached, Adaptive-filters keep readjusting the budgets due to periodic shrinking of error bounds. Finally, when 90% are stable, STAR gives more than an order of magnitude reduction in load compared to both Adaptive-filter and uniform allocation.

Next, we compare the performance of STAR, Adaptive-filters, and uniform allocation under different configurations by varying input data distribution, standard deviation (step sizes), and update frequency at each node. For data distribution, the workload is either generated from a random-walk pattern or Gaussian. For standard deviation/step-size, 70% of the nodes have uniform parameters as previously described; the remaining 30% nodes have these parameters proportional to *rank* or randomly assigned from the range [0.5, 150].

Figure 10 shows the corresponding results for different settings of data distribution and standard deviation for a 4 level degree-4 tree but fixed update frequency of 1 update per node per round. We make the following observations from Figure 10(a). First, when error budget is smaller than noise, no algorithm in any configuration achieves better performance than uniform allocation. Adaptive-filters, however, incurs a slightly higher overhead due to self-tuning even though it does not benefit. In comparison, STAR avoids self-tuning costs via cost-benefit throttling. Second, Adaptive-filters and uniform error allocation reach a cross-over point having a similar performance. This cross-over implies that for Adaptive-filters, the cost of self-tuning is equal to the benefits. Third, as error budget increases, STAR achieves better performance than Adaptive-filters. Because step-sizes are based on node rank, STAR’s outlier detection avoids allocating budget to the nodes having the largest step-sizes. Adaptive-filters, however, do not make such a distinction and computes burden scores based on load thereby favoring nodes with relatively large step sizes. Thus, since the total budget is limited, reducing error budget at nodes with small step sizes increases their load but does not benefit outliers since the additional slack is still insufficient to filter their updates. Finally, as expected, when error budget is higher than noise, all algorithms achieve better performance. In this configuration, Adaptive-filters reduces monitoring load by 3x-5x compared to uniform allocation and by 2x-3x com-

pared to Adaptive-filters.

Under random distribution of step-sizes as described above, STAR reduces load by up to 1.5x compared to Adaptive-filters and up to 3x against uniform allocation. Comparing across configurations, all algorithms perform better under input distribution of Gaussian compared to the random-walk model. Overall, across all configurations, STAR reduces monitoring load by up to an order of magnitude compared to uniform allocation and by up to 5x compared to Adaptive-filters.

5.2 Testbed experiments

In this subsection we quantify the reduction in monitoring load due to self-tuning AI and the query precision of reported results for the Distributed Heavy Hitter (DHH) monitoring application.

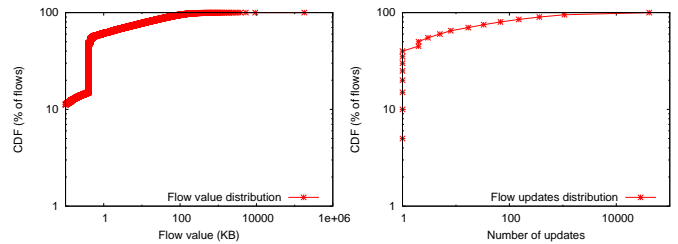


Figure 11: CDF of percentage of flows vs (a) number of updates and (b) flow-values for the Abilene dataset. The graph is on a log-log scale.

We use multiple netflow traces obtained from the Abilene [1] Internet2 backbone network. The data was 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 query on this dataset for tracking the top 100 flows (destination IP as key) in terms of bytes received over a 30 second moving window shifted every 10 seconds.

We first describe the characteristics of this Abilene workload. Figure 11 shows the cumulative distribution function (CDF) of the percentage of network flows versus the number of bytes (KB) sent by each flow. We observe that 60% flows send less than 1 KB of aggregate traffic, 80% flows send less than 12 KB traffic, 90% flows less than 55 KB and 99% of the flows send less than 330 KB of total traffic during the 1-hour run. Note that the distribution is heavy-tailed and maximum aggregate flow value is about 179.4 MB. Figure 11

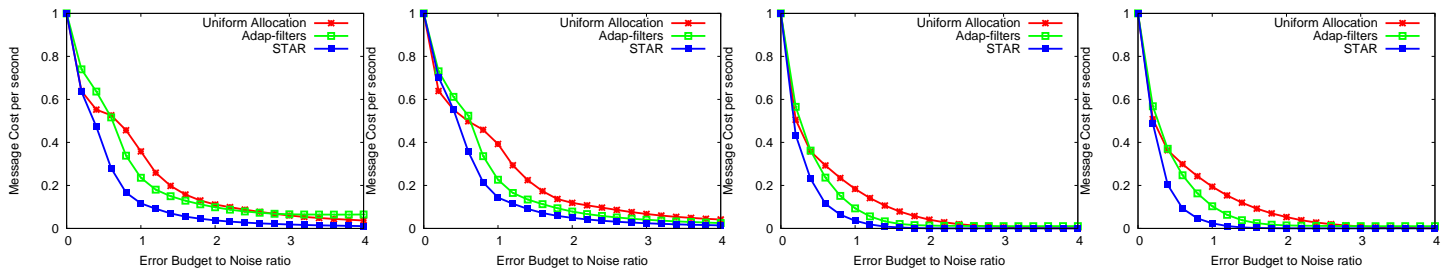


Figure 10: Performance comparison of STAR vs. Adaptive-filters and uniform allocation for different (workload, standard deviation) configurations (a) random walk, rank (b) random walk, random (c) Gaussian, rank (d) Gaussian, random.

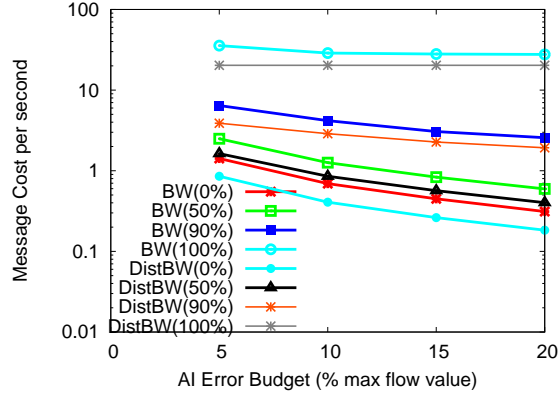


Figure 12: Performance comparison of different policies for setting initial error budgets for DHH application. BW is bandwidth cost per node and DistBW is redistribution overhead per node.

shows the corresponding CDF graph for of the percentage of network flows versus the number of updates. We observe that 40% flows send only a single update (a 28 byte IP/UDP packet). Further, 80% flows send less than 70 updates, 90% flows less than 360 updates, and 99% flows less than 2000 updates. Note that the number of update distribution is also heavy-tailed with the maximum number of updates send by a flow is about 42,000.

Overall, the total number of updates sent by roughly 80,000 flows originating from 120 sensors is around 25 million. Thus, the monitoring load for AI of 0 error budget would be about 58.6 messages per second per node for each of the 120 nodes. Therefore, a centralized scheme would incur prohibitive cost of processing this workload generating about 7,000 updates per second.

To address this scalability challenge, we analyze different settings of the optimization discussed in Section 4.3 to filter majority mice flows in our monitoring system. For this experiment, we enable distribution mechanism for error distribution using our self-tuning algorithm. Figure 12 shows the performance comparison graph by varying the $Root_{share}$ parameter values of 0%, 50%, 90%, and 100% for setting initial error budgets and global AI error budget of 5%, 10%, 15%, and 20% in each aggregation tree. Note that this figure shows both the total bandwidth cost per node (BW) and the error redistribution overhead per node (DistBW). We observe that with a $Root_{share}$ of 50% and AI of 5%, we

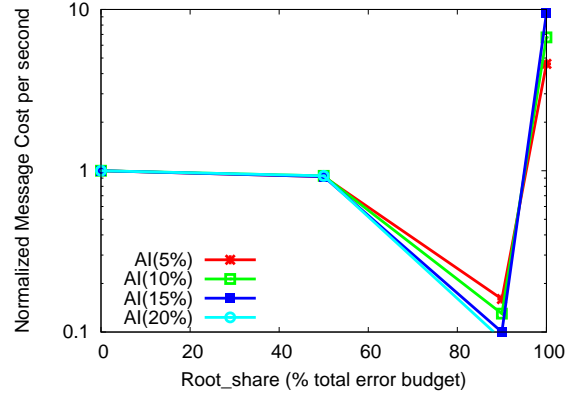


Figure 13: Performance benefits of keeping local error budget at the internal nodes of the aggregation tree.

incur overhead of about 2.5 messages per node per second. By increasing the AI error budget to 20%, we can reduce this cost by almost a factor of five to about 0.5 messages per node per second. However, for $Root_{share}$ of 100%, each mice update reaches the root which then initiates an error allocation to each node of the aggregation tree. Finally, note that even a setting of $Root_{share}$ of 90% gives an order of magnitude reduction in load compared to $Root_{share}$ of 100% by filtering mice flows at lower levels of the aggregation tree.

Next, we show the performance benefits of keeping a non-zero local error budget (δ_{self}) at the internal nodes in the tree. Intuitively, when an internal node aggregates across its children values, even though the child values have deviated significantly to bypass their own error range, the net effect by merging all children updates may still be close to zero i.e., the children updates cancel out each other. Figure 11 shows the benefits of keeping 10% error budget at the internal nodes during initial error allocation compared to 0% budget at the internal nodes. The different lines in the graph correspond to different AI budgets. We observe that for $Root_{share}$ of 50%, we achieve only 10% load reduction due to filtering at internal nodes. However, for $Root_{share}$ of 90%, we achieve almost an order of magnitude load reduction. In this case, the δ_{self} at the internal nodes is sufficient to filter a significant fraction of flows. However, when $Root_{share}$ of 50%, the local error budget is only able to filter about 10% updates. Note that this reduction due to internal nodes is in

addition to the filtering at the leaf nodes. Overall, a scalable monitoring service can significantly benefit from filtering at both the leaf and the internal nodes.

In summary, our evaluation shows that adaptive setting of modest AI budgets can provide large bandwidth savings to enable scalable monitoring.

6. RELATED WORK

Olston et al. [28] proposed a self-tuning algorithm for a *one-level* tree: increase δ for nodes with high load and low previous δ and decrease δ for nodes with low load and high previous δ . Our algorithm differs from Olston's in two fundamental ways: First, to work with large-scale multi-level trees, we define a distributed algorithm in which internal nodes not only split δ_c among their children c but may also retain δ_{self} to help prevent updates received from children from being propagated further up the tree. Second, in order to scale to tens of thousands of attributes required by per-flow network monitoring services, our algorithm performs cost/benefit throttling. In particular, we use a competitive analysis to ensure that rebalancing δ accounts for only a fraction of the system's total message load by only redistributing δ to correct large imbalances or to correct stable, long-duration (though perhaps small) imbalances.

For hierarchical topologies, Manjhi et al. [27] determine an optimal but *static* distribution of slack to the internal and leaf nodes of a tree for finding frequent items in data streams. IrisNet [12] filters sensors at leaves and caches timestamped results in a hierarchy with queries that specify the maximum staleness they will accept and that trigger retransmission if needed. Deligiannakis et al. [11] propose an adaptive precision setting technique for hierarchical aggregation, with a focus on sensor networks. However, similar to Olston's approach, their technique also periodically shrinks the error budgets for each tree which limits scalability for tracking large number of attributes. Further, since their approach uses only two *local* anchor points around the current error budget to infer the precision-performance tradeoff as shown in Figure 2, it cannot infer the complete correlation making it susceptible to dynamic changes. None of the previous studies to our knowledge uses the variance and the update rate in the data distribution in a principled manner. In comparison, STAR provides an efficient and practical algorithm that uses an informed mathematical model to compute globally optimal budgets and performs cost-benefit throttling to adaptively set precision constraints in a general communication hierarchy.

Some Recent studies [19, 22, 24] have proposed monitoring systems with distributed triggers that fire when an aggregate of remote-site behavior exceeds an a priori global threshold. Their solution is based on a centralized architecture. STAR may enhance such efforts by providing a scalable way to track top-k and other significant events.

Other studies have proposed prediction-based techniques for data stream filtering e.g., using Kalman filters [21], neural networks [25], etc. Recently, there has been a considerable interest in the database and sensor network communities on approximate data management techniques; Skordylis et al. [35] provide a good survey.

There are ongoing efforts similar to ours in the P2P and databases community to build global monitoring services. PIER is a DHT-based relational query engine [20] targeted at querying real-time data from many vantage-points on the

Internet. Sophia [38] is a distributed monitoring system designed with a declarative logic programming model. Gigascope [10] provides a stream database functionality for network monitoring applications.

Traditionally, DHT-based aggregation is event-driven and best-effort, i.e., each update event triggers re-aggregation for affected portions of the aggregation tree. Further, systems often only provide eventual consistency guarantees on its data [37, 40], i.e., updates by a live node will eventually be visible to probes by connected nodes.

7. CONCLUSIONS AND FUTURE WORK

Without adaptive setting of precision constraints, large scale network monitoring systems may be too expensive to implement (because too many events flow through the system) even under considerable error budgets. STAR provides self-tuning arithmetic imprecision to adaptively bound the numerical accuracy in query results, and optimizations to enable scalable monitoring of large number of stream events in a distributed system.

While our self-tuning algorithm focuses on minimizing communication load in an aggregation hierarchy under fixed data precision constraints, it might be useful for some applications and environments to investigate the dual problem of maximizing data precision subject to constraints on availability of global computation and communication resources. Another interesting topic of future work is to consider self-tuning error distribution for general graph topologies e.g., DAGs, rings, etc., that are more robust to node failures than tree networks. Finally, reducing monitoring load in real-world systems would also require understanding other dimensions of imprecision such as temporal where queries can tolerate bounded staleness and topological when only a subset of nodes are needed to answer a query.

8. REFERENCES

- [1] Abilene internet2 network.
<http://abilene.internet2.edu/>.
- [2] R. Avnur and J. M. Hellerstein. Eddies: continuously adaptive query processing. In *SIGMOD*, 2000.
- [3] B. Babcock, S. Babu, M. Datar, R. Motwani, and J. Widom. Models and issues in data stream systems. In *PODS*, 2002.
- [4] B. Babcock and C. Olston. Distributed top-k monitoring. In *SIGMOD*, June 2003.
- [5] S. Babu, R. Motwani, K. Munagala, I. Nishizawa, and J. Widom. Adaptive ordering of pipelined stream filters. In *SIGMOD*, 2004.
- [6] A. Bhambe, M. Agrawal, and S. Seshan. Mercury: Supporting Scalable Multi-Attribute Range Queries. In *SIGCOMM*, Portland, OR, August 2004.
- [7] S. Chaudhuri, G. Das, and V. Narasayya. A robust, optimization-based approach for approximate answering of aggregate queries. In *SIGMOD '01*, 2001.
- [8] D. D. Clark, C. Partridge, J. C. Ramming, and J. Wroclawski. A knowledge plane for the internet. In *SIGCOMM*, 2003.
- [9] R. Cox, A. Muthitacharoen, and R. T. Morris. Serving DNS using a Peer-to-Peer Lookup Service. In *IPTPS*, 2002.
- [10] C. D. Cranor, T. Johnson, O. Spatscheck, and V. Shkapenyuk. Gigascope: A stream database for

- network applications. In *SIGMOD*, 2003.
- [11] A. Deligiannakis, Y. Kotidis, and N. Roussopoulos. Hierarchical in-network data aggregation with quality guarantees. In *EDBT*, 2004.
- [12] A. Deshpande, S. Nath, P. Gibbons, and S. Seshan. Cache-and-query for wide area sensor databases. In *Proc. SIGMOD*, 2003.
- [13] C. Estan and G. Varghese. New directions in traffic measurement and accounting. In *SIGCOMM*, 2002.
- [14] M. J. Freedman and D. Mazires. Sloppy Hashing and Self-Organizing Clusters. In *IPTPS*, Berkeley, CA, February 2003.
- [15] FreePastry. <http://freepastry.rice.edu>.
- [16] Y. Fu, J. Chase, B. Chun, S. Schwab, and A. Vahdat. SHARP: An architecture for secure resource peering. In *Proc. SOSP*, Oct. 2003.
- [17] 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.
- [18] J. M. Hellerstein, V. Paxson, L. L. Peterson, T. Roscoe, S. Shenker, and D. Wetherall. The network oracle. *IEEE Data Eng. Bull.*, 28(1):3–10, 2005.
- [19] L. Huang, M. Garofalakis, A. D. Joseph, and N. Taft. Communication-efficient tracking of distributed cumulative triggers. In *ICDCS*, 2007.
- [20] R. Huebsch, J. M. Hellerstein, N. Lanham, B. T. Loo, S. Shenker, and I. Stoica. Querying the Internet with PIER. In *VLDB*, 2003.
- [21] A. Jain, E. Y. Chang, and Y.-F. Wang. Adaptive stream resource management using kalman filters. In *SIGMOD*, 2004.
- [22] A. Jain, J. M. Hellerstein, S. Ratnasamy, and D. Wetherall. A wakeup call for internet monitoring systems: The case for distributed triggers. In *HotNets*, San Diego, CA, November 2004.
- [23] N. Jain, D. Kit, P. Mahajan, P. Yalagandula, M. Dahlin, and Y. Zhang. STAR: self tuning aggregation for scalable monitoring (extended). Technical Report TR-07-15, UT Austin Department of Computer Sciences, March 2007.
- [24] R. Keralapura, G. Cormode, and J. Ramamirtham. Communication-efficient distributed monitoring of thresholded counts. In *SIGMOD*, 2006.
- [25] V. Kumar, B. F. Cooper, and S. B. Navathe. Predictive filtering: a learning-based approach to data stream filtering. In *DMSN*, 2004.
- [26] 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.
- [27] A. Manjhi, V. Shkapenyuk, K. Dhamdhere, and C. Olston. Finding (Recently) Frequent Items in Distributed Data Streams. In *ICDE*, 2005.
- [28] C. Olston, J. Jiang, and J. Widom. Adaptive filters for continuous queries over distributed data streams. In *SIGMOD*, 2003.
- [29] C. Olston and J. Widom. Offering a precision-performance tradeoff for aggregation queries over replicated data. In *VLDB*, Sept. 2000.
- [30] C. G. Plaxton, R. Rajaraman, and A. W. Richa. Accessing Nearby Copies of Replicated Objects in a Distributed Environment. In *ACM SPAA*, 1997.
- [31] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content Addressable Network. In *SIGCOMM*, 2001.
- [32] A. Rowstron and P. Druschel. Pastry: Scalable, Distributed Object Location and Routing for Large-scale Peer-to-peer Systems. In *Middleware*, 2001.
- [33] J. Shneidman, P. Pietzuch, J. Ledlie, M. Roussopoulos, M. Seltzer, and M. Welsh. Hourglass: An infrastructure for connecting sensor networks and applications. Technical Report TR-21-04, Harvard Technical Report, 2004.
- [34] A. Singla, U. Ramachandran, and J. Hodgins. Temporal notions of synchronization and consistency in Beehive. In *Proc. SPAA*, 1997.
- [35] A. Skordylis, N. Trigoni, and A. Guitton. A study of approximate data management techniques for sensor networks. In *WISES, Fourth Workshop on Intelligent Solutions in Embedded Systems*, 2006.
- [36] 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.
- [37] R. van Renesse, K. Birman, and W. Vogels. Astrolabe: A robust and scalable technology for distributed system monitoring, management, and data mining. *TOCS*, 21(2):164–206, 2003.
- [38] M. Wawrzoniak, L. Peterson, and T. Roscoe. Sophia: An Information Plane for Networked Systems. In *HotNets-II*, 2003.
- [39] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. An integrated experimental environment for distributed systems and networks. In *Proc. OSDI*, pages 255–270, Boston, MA, Dec. 2002.
- [40] P. Yalagandula and M. Dahlin. A scalable distributed information management system. In *Proc SIGCOMM*, Aug. 2004.
- [41] P. Yalagandula, P. Sharma, S. Banerjee, S.-J. Lee, and S. Basu. S³: A Scalable Sensing Service for Monitoring Large Networked Systems. In *Proceedings of the SIGCOMM Workshop on Internet Network Management*, 2006.
- [42] H. Yu and A. Vahdat. Design and evaluation of a continuous consistency model for replicated services. In *OSDI*, pages 305–318, 2000.
- [43] H. Yu and A. Vahdat. Design and evaluation of a conit-based continuous consistency model for replicated services. *ACM Trans. on Computer Systems*, 20(3), Aug. 2002.
- [44] B. Y. Zhao, J. D. Kubiatowicz, and A. D. Joseph. Tapestry: An Infrastructure for Fault-tolerant Wide-area Location and Routing. Technical Report UCB/CSD-01-1141, UC Berkeley, Apr. 2001.