

Scaling Construction of Low Fan-out Overlays for Topic-based Publish/Subscribe Systems

Chen Chen
Department of Electrical and
Computer Engineering
University of Toronto
chen@msrg.utoronto.ca

Roman Vitenberg
Department of Informatics
University of Oslo, Norway
romanvi@ifi.uio.no

Hans-Arno Jacobsen
Department of Electrical and Computer Engineering
Department of Computer Science
University of Toronto
jacobsen@eecg.toronto.edu

Abstract—It is a key challenge and fundamental problem in the design of distributed publish/subscribe systems to construct the underlying dissemination overlay. In this paper, we focus on effective practical solution for the **MinMax-TCO** problem: Create a topic-connected pub/sub overlay in which all nodes interested in the same topic are organized in a directly connected dissemination sub-overlay while keeping the maximum node degree to the minimum.

Previously known solutions provided an extensive analysis of the problem and an algorithm that achieves a logarithmic approximation for **MinMax-TCO**. Yet, they did not focus on efficiency of the solution or feasibility of decentralized operation that would not require full knowledge of the system. Compared to these solutions, our proposed algorithm produces an overlay with marginally higher degrees. At the same time, it has drastically reduced runtime cost, which is corroborated by both theoretical analysis and empirical evaluation. The latter shows a speedup by a factor of more than 25 on average for typical pub/sub workloads.

I. INTRODUCTION

We are witnessing an increasingly widespread use of the publish/subscribe (pub/sub) communication paradigm in the design of large-scale distributed systems. Pub/Sub is regarded as a technology enabler for a loosely coupled form of interaction among many publishing data sources and many subscribing data sinks. Many applications report benefits from using this form of interaction, such as application integration [1], financial data dissemination [2], RSS feed distribution and filtering [3], [4], and business process management [5]. As a result, many industry standards have adopted pub/sub as part of their interfaces. Examples of such standards include WS Notifications, WS Eventing, the OMG's Real-time Data Dissemination Service, and the Active Message Queuing Protocol.

In pub/sub, subscribers convey their interests in receiving messages and publishers disseminate publication messages. The language and data model for subscriptions and publications vary across systems. In this paper, we focus on the topic-based pub/sub model. In a topic-based system, publication messages are associated with topics and subscribers register their interests in receiving all messages published on topics of interest. Many commercial systems follow this design. For example, TIBCO RV [2] has been used extensively for market data feed dissemination and Google's GooPS [1] and Yahoo's

YMB [6] constitute the distributed message exchange for Web-based applications operating worldwide.

In a distributed pub/sub system, so called pub/sub brokers, often connected in a federated manner as an application-level overlay network, efficiently route publication messages from data sources to sinks. The distributed design was introduced to address pub/sub system scalability. The overlay of a pub/sub system directly impacts the system's performance and the message routing cost. Constructing a high-quality broker overlay is a key challenge and fundamental problem for distributed pub/sub systems that has received attention both in industry [1], [6] and academia [7], [8], [9], [10], [11], [12].

In [7], the authors defined the notion of topic connectivity, which informally speaking means that all nodes (i.e., pub/sub brokers) interested in the same topic are organized in a connected dissemination sub-overlay. This property ensures that nodes not interested in a topic would never need to contribute to disseminating information on that topic. Publication routing atop such overlays saves bandwidth and computational resources otherwise wasted on forwarding messages of no interest to the node. It also results in smaller routing tables.

An additional desirable property for a pub/sub overlay is to have a low node degree. High node degrees increase the probability of hotspots and aggravate the impact of node failures on the system. A node with a high number of adjacent links has to maintain those links (i.e., monitor the links and the neighbors [7], [9]). While overlay designs for different applications might be principally different, they all share the strive for maintaining bounded node degrees, whether in DHTs [13], wireless networks [14], or for survivable network design [15].

Unfortunately, the properties of topic-connectivity and low node degree are at odds with each other. Intuitively, a sparse overlay is unlikely to be topic-connected while a dense overlay is suboptimal with respect to the node degree. In light of this dichotomy, Onus and Richa [8] introduced the fundamental **MinMax-TCO** publish/subscribe problem: *Build a topic-connected overlay (TCO) such that the overlay degree (the maximal degree of any node in the overlay) is minimal*. Onus and Richa proved that there exists no efficient solution for overlay construction that guarantees constant-factor approximation for **MinMax-TCO** (unless $P=NP$). In face of

this impossibility result, the authors propose the MinMax-ODA algorithm and a neat proof that it achieves logarithmic approximation [8].

While the results of [8] establish a fundamental baseline for any MinMax-TCO algorithm and logarithmic approximation works sufficiently well in most cases, it is vital for a practical solution to consider a number of additional design issues. The running time of MinMax-ODA is $O(|V|^4|T|)$ wherein $|V|$ is the number of nodes and $|T|$ is the number of topics in the system. This makes the overlay construction prohibitively expensive. Furthermore, MinMax-ODA is centralized and requires complete knowledge of all the nodes in the system and their interests.

The main contribution of this paper is the design of a MinMax-TCO solution that focuses on efficiency and alleviates the above limitations. In its core lies a divide-and-conquer approach to the problem: We partition the set of nodes into subsets, build a TCO for each subset, and combine all TCOs into a global overlay. The appeal of this scheme is in the substantially faster TCO construction for each subset of nodes that requires only partial knowledge about the nodes within the partition. Since the creation of TCOs for different partitions is independent, the process can be parallelized and decentralized. Yet, in order to apply this approach we need to overcome a number of obstacles.

The first challenge is comprised in the impact of partitioning on the quality of the solution. We show that the minimal overlay degree is very sensitive to partitioning and it may increase by up to a factor of $\Theta(|T|)$ in the worst case. Our solution is based on the study of workloads in existing pub/sub systems and the observation that in practical pub/sub deployments, only a relatively small number of nodes might be interested in a large number of topics [3]. We formalize this as an assumption and optimize our solution for this case. This assumption does not simplify the MinMax-TCO problem: the number of bulk subscribers is still too large to make any brute force solution around the impossibility result effective. Furthermore, it does not reduce the running time of MinMax-ODA sufficiently for practical applications. Yet, it allows us to come up with a partitioning scheme for which the divide-and-conquer approach retains the logarithmic upper bound on the overlay degree provided by MinMax-ODA.

Next, we devise an algorithm for the combine phase. Our first solution is an adaptation of the MinMax-ODA algorithm along with the proof that the adapted algorithm preserves the approximation ratio. This solution serves as a baseline for analyzing performance, identifying bottlenecks and weaknesses, and devising more advanced algorithms. Unfortunately, it does not improve the running time of MinMax-ODA. Furthermore, it still requires global knowledge about the interests of each node.

To address this issue, we observe that not all nodes need to participate in the combine phase. In each partition, we can select a number of representative nodes so that their combined interest covers the interest of all nodes in the partition. We show that if the combine phase is only performed on the

representative nodes (one representative set from each partition), then the resulting overlay will still be topic-connected. Running the combine phase only on representative nodes drastically improves the running time and eliminates the need for a central point of control that possesses complete knowledge about the system. At the same time, it may have a profound impact on the overlay degree unless we select representative nodes in a careful and controlled way. We show how to perform this selection so as to tread the balance between overlay degree and running time.

We evaluated our solution through a series of simulations on characteristic pub/sub workloads with up to 8 000 nodes and 1 000 topics. The results indicate that on average, our solution requires less than 4.0% of the running time of known state-of-the-art algorithms while yielding an insignificant increase in the maximum node degree of 2.0. While we did not analyze space complexity or measure the program footprint, the improvements in this respect are also noteworthy: For the same pub/sub workload distribution with 8 000 nodes and under the same environmental settings, our solution was taking less than a minute to construct the overlay whereas known state-of-the-art algorithms would experience memory-related problems (with 14GB of RAM allocated).

II. RELATED WORK

Traditionally, research in the area of distributed pub/sub systems has been focusing on the efficiency and scalability of message dissemination from numerous publishers to a large number of subscribers [16], [17], [18], [19]. A more recent research direction is to consider the fundamental properties of the underlying overlays for pub/sub [7], [8], [9], [10], [11], [12], [20]. This is the direction we pursue in this paper.

Topic-connectivity is explicitly stated as a desirable property for pub/sub overlays in [7], [8], [9], [12], [21]. A number of additional topic-based pub/sub systems (e.g., [17], [22]) construct an overlay per-topic thereby attaining topic connectivity without explicitly discussing this property. A related concept of creating overlay links according to node interests has been explored in [11], [20], [23], [24], [25].

The authors of [7] introduced the Minimum Topic Connected Overlay (MinAvg-TCO) problem, which aims at constructing a topic-connected overlay with minimum number of edges, i.e., minimizing the average node degree. They proved the problem is NP-complete and presented a greedy algorithm which achieves a logarithmic approximation ratio for the average node degree. In our previous work [12], a divide-and-conquer algorithm is developed for MinAvg-TCO, which dramatically improves the running time at the expense of a minor increase in the average node degree.

The approach of this paper is different from [12] in several significant ways. The underlying MinMax-ODA algorithm exhibits principally different behavior compared to the greedy MinAvg-TCO solution of [7], which leads to different analytical results with respect to both node degree and running time. The MinMax-TCO problem itself is also different from MinAvg-TCO. In particular, we show that the maximum node

degree is much more sensitive to partitioning and other elements employed in the divide-and-conquer approach compared to the average node degree. Key elements of our solution for MinAvg-TCO such as the coverage set, produce inadequate results for MinMax-TCO. In order to address this challenge, we developed new techniques, such as the division of nodes into bulk and lightweight subscribers and a representative set with a coverage factor.

The motivation for defining MinMax-TCO in [8] is that existing algorithms for MinAvg-TCO may produce an overlay in which edges are unevenly distributed across nodes. The authors point out that the maximum node degree could be as bad as $\Theta(|V|)$ compared to that of the optimal solution. More recently, Onus and Richa [9] introduced another problem, Low-TCO, which simultaneously considers average and maximum node degree. The solution for Low-TCO proposed in [9] achieves a sub-linear approximation on both maximum and average node degrees.

III. BACKGROUND

In this section we present some definitions and background information essential for the understanding of the algorithms developed in this paper.

Let V be the set of nodes and T be the set of topics. The interest function Int is defined as $Int : V \times T \rightarrow \{true, false\}$. Since the domain of the interest function is a Cartesian product, we also refer to this function as an interest matrix. Given an interest function Int , we say that a node v is interested in some topic t if and only if $Int(v, t) = true$. We then also say that v subscribes to t . The topic set which the node v subscribes to is denoted as T_v , and we call $|T_v|$ the *subscription size* of node v .

An overlay network $G(V, E)$ is an undirected graph over the node set V with the edge set $E \subseteq V \times V$. Given an overlay network $G(V, E)$, an interest function Int , and a topic $t \in T$, we say that a subgraph $G_t(V_t, E_t)$ of G is *induced by t* if $V_t = \{v \in V | Int(v, t)\}$ and $E_t = \{(v, w) \in E | v \in V_t \wedge w \in V_t\}$. An overlay G is called *topic-connected* if for each topic $t \in T$, the subgraph G_t of G induced by t contains at most one connected component.

Beside topic-connectivity, it is important to keep the degrees of the nodes in the overlay low. Onus and Richa [8] introduced the MinMax-TCO problem for minimizing the maximum degree in a topic-connected overlay. The formal definition of the problem is as follows.

Definition 1. *MinMax-TCO(V, T, Int):* Given a set of nodes V , a set of topics T , and the interest function Int , construct a topic-connected overlay network $TCO(V, T, Int, E)$ with the smallest possible maximum node degree.

MinMax-TCO was proven to be NP-complete, and it can not be approximated by a polynomial time algorithm within a constant factor unless $P=NP$ [8]. Onus *et al.* proposed the MinMax-ODA algorithm, which always delivers a TCO that has a maximum node degree within at most $\log(|V||T|)$ times the minimum possible maximum node degree for any

TCO. Here, we refer to the MinMax-ODA algorithm by the shorter name, GM-M (Greedy algorithm for MinMax-TCO), for consistency with our notation. GM-M is specified in Algorithm 1 and operates in a greedy manner as follows: It starts with an empty set of edges and iteratively adds carefully selected edges one by one until topic-connectivity is attained. The edge selection criterion is as follows: If there exist edges whose addition to the overlay does not increase the maximum node degree, the algorithm picks an edge with the largest contribution from the set of all such edges. Otherwise, an edge with the largest contribution among all edges is selected. The contribution of an edge e , denoted as $contrib(e)$, is defined as the number of topic-connected components reduced by adding the edge to the current overlay.

Algorithm 1 Greedy algorithm for MinMax-TCO

GM-M($I(V, T, Int)$)

Input: $I(V, T, Int)$

Output: A topic-connected overlay $TCO(V, T, Int, E_{GM})$

```

1:  $E_{pot} \leftarrow \emptyset$ 
2: for all  $e = (v, w)$  s.t.  $(w, v) \notin E_{pot}$  do
3:   add  $e$  to  $E_{pot}$ 
4:  $E_{new} \leftarrow \text{buildMMEdges}(V, T, Int, \emptyset, E_{pot})$ 
5: return  $TCO(V, T, Int, E_{new})$ 

```

Algorithm 2 Overlay Construction for GM-M Algorithm

buildMMEdges($V, T, Int, E_{cur}, E_{pot}$)

Input: $V, T, Int, E_{cur}, E_{pot}$

// E_{cur} : Set of current edges that exist in the overlay

// E_{pot} : Set of potential edges that can be added

Output: Edge set E_{new} that combined with E_{cur} , forms a TCO

```

1:  $E_{new} \leftarrow \emptyset$ 
2: for all  $e = (v, w) \in E_{pot}$  do
3:    $contrib(e) \leftarrow |\{t \in T | Int(v, t) \wedge Int(w, t) \wedge$ 
    $v, w \text{ belong to different connected components for } t \text{ in } G(V, E_{cur})\}|$ 
4: while  $G(V, E_{new} \cup E_{cur})$  is not topic-connected do
5:    $e \leftarrow$  find edge  $e$  s.t.  $contrib(e)$  is maximum and  $e$  increases the
   maximum degree of  $G(V, E_{new} \cup E_{cur})$  minimally
6:    $E_{new} \leftarrow E_{new} \cup \{e\}$ 
7:    $E_{pot} \leftarrow E_{pot} - \{e\}$ 
8: for all  $e = (v, w) \in E_{pot}$  do
9:    $contrib(e) \leftarrow$  update the contribution of a potential edge  $e$  as
   the reduction on the number of topic-connected components which
   would result from the addition of  $e$  to  $G(V, E_{new} \cup E_{cur})$ 
10: return  $E_{new}$ 

```

The following results about GM-M were proven in an elegant fashion in [8]:

Lemma 1 (GM-M Approximation Ratio & Running Time). *The overlay network output by Algorithm 1 has a maximum node degree within a factor of $\log(|V||T|)$ of the maximum node degree of the optimal solution for MinMax-TCO(V, T, Int). Algorithm 1 has a running time of $O(|E_{pot}|^2|T|) = O(|V|^4|T|)$.*

IV. DIVIDE-AND-CONQUER FOR MINMAX-TCO

Taking into account the GM-M running time of $O(|V|^4|T|)$, the number of nodes is the most significant factor determining the performance and scalability of the solution for MinMax-TCO. In view of this, we devise a divide-and-conquer strategy

to solve the problem: (1) *divide* the MinMax-TCO problem into several sub-overlay construction problems that are similar to the original overlay but with a smaller node set, (2) *conquer* the sub-MinMax-TCO problems independently and build sub-overlays into sub-TCOs, and then (3) *combine* these sub-TCOs to one TCO as a solution to the original problem.

In this section we present the design steps and key decisions of the divide-and-conquer approach for MinMax-TCO. We show analytically that the resulting algorithm leads to a significantly improved running time cost and reduced knowledge requirement as compared to the GM-M algorithm. In Section VI, we quantify these improvements empirically.

A. GM-M as a Building Block for Divide-and-Conquer

In this section, we analyze the GM-M algorithm in greater depth and derive several new results about its running time. GM-M is employed as a building block in our divide-and-conquer algorithms and it serves as the baseline for our experimentation.

First, we show that the maximum node degree of the overlay produced by GM-M is bounded by the maximum subscription size in the input.

Lemma 2 (Bound on the maximum degree). *The maximum node degree of the TCO produced by GM-M is $O(\max_{v \in V} |T_v|)$.*

In [8], Algorithm 1 is the only entry point for Algorithm 2. This means that E_{cur} is always equal to \emptyset and E_{pot} to $V \times V$ upon the invocation of Algorithm 2. When we adapt GM-M for the combine phase of the divide-and-conquer approach, we need to apply GM-M on a collection of TCOs already produced in the conquer phase. Therefore, we have to extend the analysis of GM-M for the case when E_{cur} is non-empty and E_{pot} is equal to $(V \times V) \setminus E_{cur}$. Let E_{new} be the set of edges returned by Algorithm 2. Denote the maximal degree of $G(V, E)$ by $D(V, E)$ and the maximal degree of the optimal solution for MinMax-TCO(V, T, Int) by $D_{OPT}(V, T, Int)$. Then, the following result holds:

Lemma 3. *If invoked on V, T, Int, E_{cur} , and E_{pot} such that $E_{cur} \cup E_{pot} = V \times V$, Algorithm 2 outputs E_{new} such that*

- (a) $G(V, E_{cur} \cup E_{new})$ is topic-connected and
- (b) the maximum node degree $D(V, E_{cur} \cup E_{new})$ is bounded by $O(D(V, E_{cur}) + D_{OPT}(V, T, Int) \cdot \log(|V||T|))$.

We provide detailed proofs of Lemma 2 and 3 in the full version of this paper [26].

B. Divide and Conquer Phases of the Solution

There exist two principal methods to *divide* the nodes: (1) node clustering and (2) random partitioning. Node clustering is organizing the original node set into groups so that nodes with similar interests are placed in the same group while nodes with diverging interests belong to different groups. Random partitioning assigns each node in the given node set to one of the partitions based on a uniformly random distribution. Once the partitions are determined, the existing GM-M algorithm

can be employed to *conquer* the sub-MinMax-TCO problems by determining inner edges used for the construction of the sub-overlays.

The idea of node clustering seems attractive because well-clustered nodes with strongly correlated interests would result in lower maximum node degrees in the sub-TCOs produced by GM-M. The problem with this approach is the high runtime cost of clustering algorithms taking into account the large number of nodes and varying subscription size. Additionally, they require the computation of a “distance” metric among nodes. In our case this translates to calculating pairwise correlations among node interests with significant run time cost implications. It is challenging to fit node clustering into the divide-and-conquer approach so that the latter is still superior to the GM-M algorithm in terms of running time cost. Furthermore, it is difficult to devise an effective decentralized algorithm for node clustering that would not require complete knowledge of V and Int . Finally, node clustering by interests may yield clusters that vary in size depending on the clustering algorithm used. On the other hand, the divide-and-conquer approach performs optimally when partitions are equal-sized and there are no large clusters that stand out.

Algorithm 3 Naive algorithm for *divide* and *conquer* phases

Input: V, T, Int, p
// p : the number of partitions, $1 \leq p \leq |V|$.
Output: A list of TCOs $List_{TCO}$, one TCO for each partition

- 1: $List_{TCO} \leftarrow \emptyset$
- 2: Randomly divide V into p partitions $V_d, d=1, 2, \dots, p$
- 3: **for** $d = 1$ to p **do**
- 4: $Int_d \leftarrow Int|_{V_d}$
- 5: $TCO_d(V_d, T, Int_d, E_d) \leftarrow \text{GM-M}(V_d, T, Int_d)$
- 6: add TCO_d to $List_{TCO}$
- 7: **return** $List_{TCO}$

We choose random partitioning for the divide-and-conquer approach because it is extremely fast, more robust than node clustering, easier to tune, and it can be realized in a decentralized manner. Furthermore, the construction of inner edges for each overlay only requires knowledge of node interests within the overlay. Hence, random partitioning can be oblivious to the composition of nodes and their interests. The number of partitions p is given as an input parameter and each sub-overlay has $k = |V_d| = \frac{|V|}{p}$ nodes, where $d = 1, \dots, p$. This equal-sized division is optimal with respect to the running time. The resulting algorithm for the *divide* and *conquer* phases is presented in Algorithm 3.

Unfortunately, random partitioning may place nodes with diverging interest into the same partition thereby reducing the amount of correlation that is present in the original node set. As Lemma 4 shows, this may have a profound effect on the maximum node degree. Consider the overlay $G(V, E)$ for the *conquer* phase produced by this algorithm where $List_{TCO}[d] = TCO_d(V_d, T, Int_d, E_d), E = \cup_{d=1}^p E_d$. Then,

Lemma 4. *There is an instance $I(V, T, Int)$ of MinMax-TCO on which the maximum degree $D(V, E)$ of the overlay output by Algorithm 3 is greater by a factor of $\Theta(|T|)$ than*

the maximum node degree $D_{OPT}(V, T, Int)$ of the optimal solution for $I(V, T, Int)$.

Proof: Construct an instance $I(V, T, Int)$ of MinMax-TCO as follows: Consider the topic set $T = \{t_1, t_2, \dots, t_m\}$ of size $m = 2^h$. Node set V consists of $h+1 = \log m + 1$ subsets, denoted as $A_i, 0 \leq i \leq h$; each node subset A_i contains 2^i nodes, i.e., $A_i = \{v_{(i,1)}, \dots, v_{(i,2^i)}\}$. Each node $v_{(i,j)}, 1 \leq j \leq 2^i$ subscribes to a topic set $T_{(i,j)}$ of size $\frac{m}{2^i}$, defined as follows:

$$\begin{aligned}
 V &= \bigcup_{i=0}^h A_i & T &= \{t_1, \dots, t_m\} \\
 A_0 &= \{v_{(0,1)}\} & T_{(0,1)} &= \{t_1, \dots, t_m\} \\
 A_1 &= \{v_{(1,1)}, v_{(1,2)}\} & T_{(1,1)} &= \{t_1, \dots, t_{\frac{m}{2}}\}, T_{(1,2)} = \{t_{\frac{m}{2}+1}, \dots, t_m\} \\
 &\vdots & & \\
 A_i &= \{v_{(i,1)}, \dots, v_{(i,2^i)}\} & T_{(i,j)} &= \{t_{\frac{j}{2^i}m+1}, \dots, t_{\frac{(j+1)m}{2^i}}\}, 1 \leq j \leq 2^i \\
 &\vdots & & \\
 A_h &= \{v_{(h,1)}, \dots, v_{(h,m)}\} & T_{(h,1)} &= \{t_1\}, \dots, T_{(h,m)} = \{t_m\}
 \end{aligned}$$

The optimal overlay TCO_{OPT} for $I(V, T, Int)$ is shown in Fig.1(a). Its maximum node degree is a constant: $D_{OPT} = 3$. Consider the output of Algorithm 3. Due to the random partitioning in Line 2, there is a chance for generating a partition that consists of nodes in $A_0 = \{v_{(0,1)}\}$ and $A_h = \{v_{(h,1)}, \dots, v_{(h,m)}\}$ (see Fig.1(b)). To attain topic-connectivity for this partition, $v_{(0,1)}$ has to be linked to all m nodes in A_h , which makes the node degree of $v_{(0,1)}$ to be $m = |T|$.

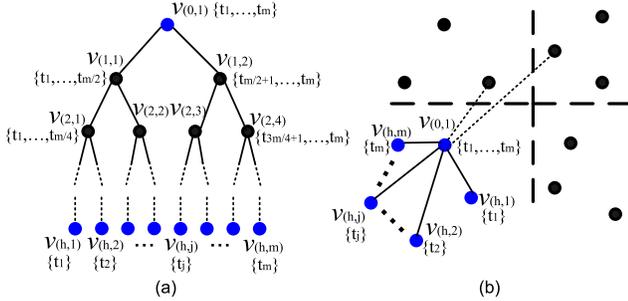


Fig. 1. (a) TCO_{OPT} with $D_{OPT} = 3$. (b) TCO for random partitioning.

Consequently, the maximum degree $D(V, E)$ of the overlay output by Algorithm 3 for $I(V, T, Int)$ is greater by a factor of $\Theta(\frac{m}{3}) = \Theta(|T|)$ than the maximum node degree $D_{OPT}(V, T, Int)$ of the optimal solution for $I(V, T, Int)$. ■

Essentially, Lemma 4 shows that if we use random partitioning for the divide-and-conquer approach, then the overlay degree for the *conquer* phase alone may exceed the overlay degree for the complete optimal solution by a factor of $\Theta(|T|)$. Furthermore, our empirical validation indicates that not only for a manually constructed worst case but also for the typical pub/sub workloads, random partitioning causes significant increase in the maximum node degree. Fig. 2 illustrates this effect for the default workload defined and motivated in Section VI. It compares the maximum degree for the *conquer* phase produced by Algorithm 3 with the total degree of the overlay produced by GM-M.

This effect of increased maximum degree occurs when a node subscribed to a large number of topics (i.e., a *bulk*

subscriber) is placed into the same partition with nodes whose subscriptions are not correlated. Then, such nodes do not benefit from creating a link to each other and need to connect to the bulk subscriber. Our solution to this problem is based on the study of representative pub/sub workloads used in actual applications that are described and characterized in [21]. According to these characterizations, the ‘‘Pareto 80–20’’ rule works for pub/sub workloads: Most nodes subscribe to a relatively small number of topics. The phenomenon of increased maximum degree due to partitioning still exists in such workloads as the example of Lemma 4 indicates. Yet, this observation allows us to devise an effective solution: We provide an algorithm that applies random partitioning only to such *lightweight* subscribers, performs the *conquer* phase for lightweight partitions, and merges them with bulk subscribers at the *combine* phase.

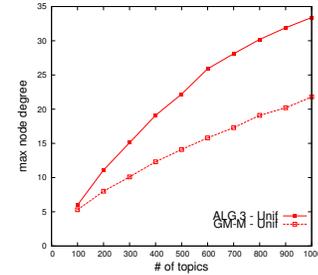


Fig. 2. Poor performance of naive divide-and-conquer

Formally, given an instance $I(V, T, Int)$ for MinMax-TCO, we introduce an additional parameter called bulk subscription threshold η , $\eta \in (0, |T|]$. η determines the division of V into the set of bulk and lightweight subscribers B and L , respectively: $B = \{v : |T_v| > \eta\}$ and $L = \{v : |T_v| \leq \eta\}$. Algorithm 4 applies random partitioning to L , creates a TCO for each of the partitions, and returns a list of these TCOs.

Algorithm 4 Divide and conquer phases for lightweight nodes

```

conquerLightweight( $I(V, T, Int), \eta, p$ )
Input:  $I(V, T, Int), \eta$ 
      //  $\eta$ : the bulk node threshold;
      //  $p$ : the number of partitions for lightweight nodes.
Output: A list of TCOs  $List_{TCO}$ , one TCO for each partition
1:  $B \leftarrow \{v \in V : |T_v| > \eta\}, Int_B \leftarrow Int|_B$ 
2:  $List_{TCO} \leftarrow \emptyset$ 
3: Randomly divide  $L = V - B$  into  $p$  partitions  $L_d, d = 1, \dots, p$ 
4: for  $d = 1$  to  $p$  do
5:    $Int_d \leftarrow Int|_{L_d}$ 
6:    $TCO_d(V_d, T, Int_d, E_d) \leftarrow \text{GM-M}(L_d, T, Int_d)$ 
7:   add  $TCO_d$  to  $List_{TCO}$ 
8: return  $List_{TCO}$ 

```

Lemma 2 can be directly applied to the overlay output by Algorithm 4 to produce a bound on its maximum degree.

Lemma 5 (Bound on the maximum degree for the *conquer* phase). *The maximum node degree of an overlay $G(V, E)$ produced by Algorithm 4 is bounded by the bulk subscription threshold: $D(V, E) = O(\eta)$.*

We now analyze the running time of Algorithm 4.

Lemma 6 (Running time for the *conquer* phase). *The running time cost of Algorithm 4 is $O(\frac{|L|^4|T|}{p^3})$.*

Proof: The loop in Lines 4–7 of Algorithm 4 uses GM-M to build a sub-TCO for each sub-overlay. Each sub-overlay has at most $\frac{|L|}{p}$ nodes so that by Lemma 1, the running time for constructing each sub-TCO is $O(\frac{|L|^4|T|}{p^4})$. Thus, the running time for constructing all p overlays is $O(\frac{|L|^4|T|}{p^3})$. ■

This algorithm is used as a building block for the complete divide-and-conquer solution to MinMax-TCO presented in Section IV-C.

C. Combine Phase of the Solution

Algorithm 5 Divide-and-Conquer with Bulk Nodes and Lightweight Nodes for MinMax

DCB-M($I(V, T, Int), \eta, p$)

Input: $I(V, T, Int), \eta$

// η : the bulk node threshold;

// p : the number of partitions for lightweight nodes.

Output: A topic-connected overlay $TCO(V, T, Int, E_{DCB})$

- 1: $List_{TCO} \leftarrow \text{conquerLightweight}(I(V, T, Int), \eta, p)$
- 2: $TCO(V, T, Int, E_{DCB}) \leftarrow \text{combineB\&L}(V, T, Int, List_{TCO})$
- 3: return $TCO(V, T, Int, E_{DCB})$

Algorithm 6 Combine B nodes and L nodes greedily

combineB\&L($V, T, Int, List_{TCO}$)

Input: $V, T, Int, List_{TCO}$

/* $List_{TCO}$: a list of p node-disjoint TCOs for lightweight nodes:

$List_{TCO}[d] = TCO_d(L_d, T, Int_d, E_d), d=1, \dots, p$ */

Output: A topic-connected overlay $TCO(V, T, Int, E_{DCB})$

- 1: $B \leftarrow V - \bigcup_{d=1}^p L_d$ // L_d is short for $List_{TCO}[d].L_d$
- 2: $E_{inDCB} \leftarrow \bigcup_{d=1}^p E_d$ // E_d is short for $List_{TCO}[d].E_d$
- 3: $E_{potDCB} \leftarrow \{e = (v, w) | (v \in B, w \in V) \wedge (w, v) \notin E_{potDCB}\}$
- 4: $E_{potDCB} \leftarrow E_{potDCB} \cup \{e = (v, w) | v \in L_i, w \in L_j, i < j\}$
- 5: $E_{outDCB} \leftarrow \text{buildMMEEdges}(V, T, Int, E_{inDCB}, E_{potDCB})$
- 6: $E_{DCB} \leftarrow E_{inDCB} \cup E_{outDCB}$
- 7: return $TCO(V, T, Int, E_{DCB})$

In the core of our design for the *combine* phase solution lies the observation that Algorithm 2 can be used to merge the sub-overlays for different partitions and the set of bulk subscribers into a single TCO. To implement this idea, we devise Algorithm 6 that applies Algorithm 2 on a union of the sub-overlays produced at the *conquer* phase by Algorithm 4. Algorithm 5 called DCB-M, presents a complete divide-and-conquer solution for MinMax-TCO.

Lemma 7. (Correctness) *Algorithm 5 is correct: it yields an overlay such that for every topic t , all nodes interested in t are organized in a single connected component.*

Given an instance $I(V, T, Int)$ for MinMax-TCO, let TCO_{DCB} be the TCO produced by Algorithm 5. Denote its maximum node degree as D_{DCB} . There are two types of edges that form the TCO_{DCB} : (1) E_{inDCB} , the inner edges constructed by Algorithm 4, $E_{inDCB} = \bigcup_{d=1}^p E_d$ and (2) E_{outDCB} , the outer edges conjoining bulk subscribers and lightweight node sub-TCOs, which are created in Line 5 of Algorithm 6. The maximum node degree induced by E_{inDCB} and E_{outDCB} are denoted as D_{inDCB} and D_{outDCB} , respectively. It holds that

$$D_{DCB} \leq D_{inDCB} + D_{outDCB}. \quad (1)$$

Equation 1 along with Lemma 5 and Lemma 3 allow us to establish an upper bound on the degree of the overlay produced by DCB-M.

Lemma 8 (Degree bound for DCB-M). *The overlay network output by Algorithm 5 has maximum node degree $D_{DCB} = O(\eta + D_{OPT}(V, T, Int) \cdot \log(|V||T|))$.*

Corollary 1 (Approximation ratio for DCB-M). *If we regard the bulk node threshold as a constant factor or if the maximum degree E_{inDCB} of the sub-overlays constructed at the conquer phase is smaller than the maximum degree $D_{OPT}(V, T, Int)$ of the optimal overlay for MinMax-TCO(V, T, Int), then*

$$D_{DCB} = D_{OPT}(V, T, Int) \cdot O(\log(|V||T|)). \quad (2)$$

If the conditions in Corollary 1 hold, then the DCB-M algorithm achieves the same logarithmic approximation ratio as the GM-M algorithm (Algorithm 1).

Next, we consider the running time of the DCB-M algorithm. Let \mathbb{T}_{inDCB} and \mathbb{T}_{outDCB} denote the running time to build E_{inDCB} and E_{outDCB} , respectively. Let \mathbb{T}_{DCB} be the total running time cost of Algorithm 5. Then, we obtain the running time of DCB-M with:

Lemma 9 (Running time for DCB-M). *The running time of Algorithm 5 is $\mathbb{T}_{DCB} = O(\mathbb{T}_{inDCB} + \mathbb{T}_{outDCB}) = O(|T| \cdot (|B||V| + |L|^2)^2)$.*

Proof: $\mathbb{T}_{inDCB} = O(\frac{|L|^4|T|}{p^3})$ following Lemma 6.

\mathbb{T}_{outDCB} is determined by Algorithm 6, whose running time is dominated by the invocation of Algorithm 2 in Line 5. Based on Lemma 1, we have:

$$\begin{aligned} \mathbb{T}_{outDCB} &= O(|T| \cdot |E_{potDCB}|^2) \\ &= O(|T| \cdot (|B||V| + |L|^2)^2) \end{aligned} \quad (3)$$

$$\begin{aligned} \mathbb{T}_{DCB} &= O(\mathbb{T}_{inDCB} + \mathbb{T}_{outDCB}) \\ &= O(\mathbb{T}_{outDCB}) = O(|T| \cdot (|B||V| + |L|^2)^2) \end{aligned} \quad (4)$$

In summary, Algorithm 5 has asymptotic performance very similar to that of Algorithm 1, both with respect to the maximum degree and running time. As both the above analysis and experimental evaluations in Section VI indicate, it produces a marginally higher overlay degree at marginally better runtime cost. Furthermore, the *combine* phase still requires complete knowledge of V and Int , which makes decentralization infeasible. These shortcomings motivate the development of an improved solution for the *combine* phase. Our improvement is based on the notion of *Representative Set*, which we explain next.

Given $I(V, T, Int)$ for MinMax-TCO and a topic $t, t \in T$, we denote the set of subscribers to t by $subs(t)$, $subs(t) = \{v | v \in V \wedge Int(v, t)\}$. Then, the notion of a *representative set* (*rep-set*) is defined as follows:

Definition 2 (Representative set). *Given $I(V, T, Int)$, a rep-set with the coverage factor λ , denoted as $R(\lambda)$ (or R), is a subset of V such that*

$$|subs(t)|_R \geq \min\{\lambda, |subs(t)|\}, \forall t \in T \quad (5)$$

A node $r \in R(\lambda)$ is referred to as a representative node (rep-node).

As illustrated in Figure 3, a rep-set is a subset of overlay nodes that represents the interests of all the nodes in the overlay. Each topic of interest is covered by at least λ subscribers in R unless the total number of subscribers to this topic is smaller than λ . The complete node set V is always a rep-set, but there might exist many other rep-sets with much fewer rep-nodes. In essence, these nodes can function as bridges for the purpose of determining cross-TCO connections. Observe that it is possible to attain full topic-connectivity only by using cross-TCO links among rep-nodes for different partitions. Suppose we have a number of TCOs, and each TCO is represented by a rep-set (of a smaller size). To achieve topic-connectivity for a topic $t \in T$, we can just connect nodes from different rep-sets which are interested in t .

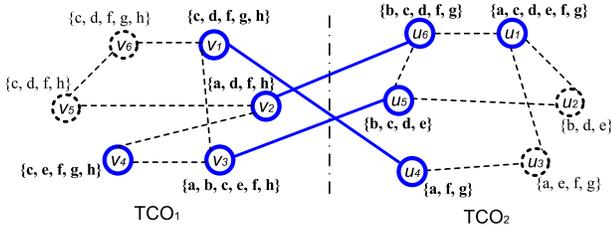


Fig. 3. $R_1 = \{v_1, v_2, v_3, v_4\}$ and $R_2 = \{u_1, u_4, u_5, u_6\}$ are rep-sets with $\lambda = 2$ for TCO_1 and TCO_2 respectively; a complete TCO for all nodes is obtained by adding cross-TCO links between R_1 and R_2 .

For typical pub/sub workloads and sufficiently large partitions, minimal rep-sets are several times smaller than the total number of nodes. This leads to significant benefits if we consider only rep-nodes as candidates for cross-TCO links. One, the running time of the overlay construction algorithms discussed in this paper is roughly proportional to the number of nodes up to the fourth degree, therefore our algorithm that only considers rep-nodes runs much faster. Two, calculation of cross-TCO links no longer requires complete knowledge of V and Int , and only a partial view of rep-nodes from rep-sets and their interests is needed. Three, rep-sets of different TCOs can be computed in parallel in a fully decentralized fashion.

At the same time, minimality of rep-sets also has an adverse effect on the maximum overlay degree due to reduced correlation across rep-sets for different TCOs. Revisit the instance $I(V, T, Int)$ of MinMax-TCO described in the example of Lemma 4. Suppose V is divided into two partitions: a partition of bulk subscribers that includes node subsets $A_i, 0 \leq i \leq h-1$ and a partition of lightweight subscribers that includes node subset A_h . The minimal rep-set for the first partition contains a single node $v_{(0,1)}$ whereas the minimal rep-set for the second partition contains all nodes in A_h . If we merge the rep-sets at the combine phase, the degree of $v_{(0,1)}$ will be $|A_h| = \Theta(|T|)$. On the other hand, if we merge the whole partitions, the optimal overlay will be the one depicted in Fig. 1(a) with constant maximum degree.

The difference arises due to the fact that in the former case, $v_{(0,1)}$ serves the focal point for all cross-overlay links while

in the latter case, these links are evenly distributed across all the nodes of the first partition. We employ two techniques to prevent the above effect. First, we only use rep-sets for the partitions of lightweight nodes and not for bulk subscribers. This is because the degree of lightweight nodes is bounded by $O(\eta)$ as we later show in Lemma 11 so that the effect is not as significant for lightweight nodes compared to bulk subscribers. Second, we use a coverage factor greater than one to ensure that there are always multiple choices when connecting the nodes for any topic.

We still need to consider, how to efficiently determine a minimal rep-set given V , T , Int , and λ . The problem of computing a minimal rep-set set is equivalent (through a linear reduction) to a variation of the classic NP-complete *Set Cover* problem, in which each item has to be covered by at least λ sets. Algorithm 9 provides a greedy implementation that attains a provable logarithmic approximation [27]. The algorithm starts with an empty rep-set and continues adding nodes to the rep-set one by one until all topics of interest are λ -covered, i.e., covered by at least λ nodes. At each iteration, the algorithm selects a node that is interested in the largest number of topics that are not yet λ -covered.

Algorithm 8 presents the resulting implementation for combining sub-TCOs. The algorithm operates in two phases. First, it determines a rep-set for each sub-TCO. Note that the rep-set for TCO_d does not need to cover all of T_d . It suffices to cover $T_{out_d} = T_d \cap (\bigcup_{i \neq d} T_i)$. Second, the algorithm connects all the nodes in the rep-sets as well as bulk subscribers into a TCO in a greedy manner by using Algorithm 2. Algorithm 7 called DCBR-M, presents our complete solution for MinMax-TCO.

Below, we establish correctness, approximation ratio and running time properties for the DCBR-M algorithm.

Lemma 10 (Correctness). *Algorithm 7 is correct: it yields an overlay such that for every topic t , all nodes interested in t are organized in a single connected component.*

Following the notations for DCB-M algorithm, we denote the TCO produced by Algorithm 7 as TCO_{DCBR} and its maximum node degree as D_{DCBR} . Observe that by operating on a reduced set of nodes at the *combine* phase, the invocation of Algorithm 2 in line 12 of Algorithm 8 solves an instance of $MinMax\text{-}TCO(BR, T, Int|_{BR})$ where BR is a union of B and all rep-sets $\bigcup_{d=1}^p R_d$. The fact that $D_{DCBR} \leq D_{inDCBR} + D_{outDCBR}$ along with Lemma 5 and Lemma 3 allow us to establish an upper bound on the degree of the overlay produced by DCBR-M.

Lemma 11 (Degree bound for DCBR-M). *The overlay network TCO_{DCBR} output by Algorithm 7 has maximum node degree: $D_{DCBR} = O(\eta + D_{OPT}(BR, T, Int|_{BR}) \cdot \log(|BR||T|))$.*

According to Lemma 11, if we choose a sufficiently large coverage factor λ so that $D_{OPT}(BR, T, Int|_{BR}) \approx D_{OPT}(V, T, Int)$, then Algorithm 7 will generate a TCO whose maximum node degree is asymptotically the same as that of the TCO output by Algorithm 5.

Algorithm 7 Divide-and-Conquer with Bulk Nodes and Lightweight Rep-nodes for MinMax

DCBR-M($I(V, T, Int), \eta, p, \lambda$)

Input: $I(V, T, Int), \eta, p, \lambda$
 // λ : the coverage factor.

Output: A topic-connected overlay $TCO(V, T, Int, E_{DCBR})$

- 1: $L_{TCO} \leftarrow \text{conquerLightweight}(I(V, T, Int), \eta, p)$
 - 2: $TCO(V, T, Int, E_{DCBR}) \leftarrow \text{combineB\&LReps}(V, T, Int, L_{TCO}, \lambda)$
 - 3: return $TCO(V, T, Int, E_{DCBR})$
-

Algorithm 8 Combine B nodes and L rep-nodes greedily

combineB\&LReps($V, T, Int, List_{TCO}, \lambda$)

Input: $V, T, Int, List_{TCO}, \lambda$
Output: A topic-connected overlay $TCO(V, T, Int, E_{DCBR})$

- 1: $B \leftarrow V - \bigcup_{d=1}^p L_d$ // L_d is short for $List_{TCO}[d].L_d$
 - 2: $E_{inDCBR} \leftarrow \bigcup_{d=1}^p E_d$ // E_d is short for $List_{TCO}[d].E_d$
 - 3: **for** $d = 1$ to p **do**
 - 4: $T_d \leftarrow \bigcup_{v \in L_d} T_v$
 - 5: **for** $d = 1$ to p **do**
 - 6: $T_{out_d} \leftarrow T_d \cap (\bigcup_{i \neq d} T_i)$
 - 7: $R_d \leftarrow \text{getRepSetFromNodes}(L_d, T_{out_d}, Int, \lambda)$
 - 8: $R \leftarrow \bigcup_{d=1}^p R_d$
 - 9: $BR \leftarrow B \cup R$
 - 10: $E_{potDCBR} \leftarrow \{e = (v, w) | (v \in B, w \in BR) \wedge (w, v) \notin E_{potDCBR}\}$
 - 11: $E_{potDCBR} \leftarrow E_{potDCBR} \cup \{e = (v, w) | v \in R_i, w \in R_j, i < j\}$
 - 12: $E_{outDCBR} \leftarrow \text{buildMMEEdges}(BR, T, Int|_{BR}, E_{inDCBR}|_{BR}, E_{potDCBR})$
 - 13: $E_{DCBR} \leftarrow E_{inDCBR} \cup E_{outDCBR}$
 - 14: return $TCO(V, T, Int, E_{DCBR})$
-

Algorithm 9 Determine a representative set for a partition

getRepSetFromNodes($V_d, T_{out_d}, Int, \lambda$)

Input: $V_d, T_{out_d}, Int, \lambda$
Output: R_d : A representative set for V_d

- 1: Start with $T_{toCover} = T_{out_d}$ and $R_d = \emptyset$
 - 2: **for all** $t \in T_{toCover}$ **do**
 - 3: $N_{toCover}[t] = \lambda$
 - 4: **while** $T_{toCover} \neq \emptyset$ **do**
 - 5: $r \leftarrow \arg \min_{v \in V_d - R_d} (\frac{1}{|\{t | t \in T_{toCover} \wedge Int(v, t)\}|})$
 - 6: $R_d \leftarrow R_d \cup \{r\}$
 - 7: **for all** $t \in T_{toCover} \wedge Int(r, t)$ **do**
 - 8: $N_{toCover}[t] \leftarrow N_{toCover}[t] - 1$
 - 9: **if** $N_{toCover}[t] = 0$ **then**
 - 10: $T_{toCover} \leftarrow T_{toCover} - \{t\}$
 - 11: Return R_d
-

Using the same reasoning as in Lemma 6, we can derive the running time cost of DCBR-M.

Lemma 12 (Running time for DCBR-M). *The running time of Algorithm 7 is $\mathbb{T}_{DCBR} = O(|T| \cdot (|B| + |R|)^4 + \frac{|L|^4}{p^3})$.*

Proof:

$$\begin{aligned} \mathbb{T}_{outDCBR} &= O(|T| \cdot |E_{potDCBR}|^2) \\ &= O(|T| \cdot (|B| \cdot (|B| + |R|) + |R|^2)^2) \\ &= O(|T| \cdot (|B| + |R|)^4) \end{aligned} \quad (6)$$

$$\begin{aligned} \mathbb{T}_{DCBR} &= O(\mathbb{T}_{outDCBR} + \mathbb{T}_{inDCBR}) \\ &= O(|T| \cdot (|B| + |R|)^4 + \frac{|L|^4}{p^3}) \end{aligned} \quad (7)$$

Lemma 12 shows that if representative sets are significantly smaller than partitions, the bulk subscribers threshold is selected so that there are few bulk subscribers, and the number

of partitions is sufficiently large, then Algorithm 7 achieves significant speedup compared to Algorithm 5. This is also corroborated by our experiments in Section VI.

D. Decentralizing the DCBR-M algorithm

Note that the DCBR-M algorithm as presented above is fully centralized. It is possible to decentralize it in the following way: (1) each lightweight node autonomously decides which random partition it belongs to and registers itself under the partition name (it is possible, e.g., to use a DHT for that purpose), (2) nodes from the same partition learn about each other and establish a communication channel, (3) different partitions construct sub-TCOs in parallel, i.e., the nodes within each partition exchange their interests and execute the GM-M algorithm, (4) different partitions compute rep-sets in parallel, (5) bulk subscribers and rep-nodes from different rep-sets communicate their interests and compute outer edges. Note that the original GM-M algorithm does not lend itself to such decentralization.

This decentralization scheme has several important benefits: reducing the total runtime cost, optimizing distributed resource utilization, and balancing the computational load. The time \mathbb{T}_{inDCBR} for computing inner edges becomes $O(\frac{|L|^4}{p^4})$ and the total time \mathbb{T}_{DCBR} becomes

$$\mathbb{T}_{DCBR} = O(|T| \cdot (|B| + |R|)^4 + \frac{|L|^4}{p^4}). \quad (8)$$

Furthermore, decentralization eliminates the need for a central entity that must maintain a global knowledge of all nodes and their interests. To quantify this benefit, we introduce additional performance characteristic of the algorithm called *potential neighbor set*, which is the set of other nodes a node has to learn about in the course of the algorithm. This characteristic is important because gathering nodes' interests in a scalable and robust decentralized manner is a problem in its own right. Additionally, the fan-out of node v in the overlay produced by *any* algorithm cannot exceed the size of the potential neighbor set of v . Therefore, minimizing the potential neighbor set has an additional desirable effect from this point of view.

To formalize this argument, we define the *potential neighbor ratio* for a node v , denoted as $pn\text{-ratio}(v)$. Potential neighbor ratio is the size of potential neighbor set for v (including v itself) normalized by the total number of nodes $|V|$. For any centralized algorithm, this ratio is equal to 1. For DCBR-M, the potential neighbor set for v consists of three subsets: (1) nodes in the same partition as v : V_d such that $v \in V_d$ (if v is a lightweight node); (2) bulk subscribers B (if v is a bulk subscriber itself or it belongs to some rep-set); and (3) all rep-nodes from other partitions: $\{u | u \in R_i \text{ s.t. } v \in B \vee v \in R_d \wedge R_i \neq R_d\}$ (if v is a bulk subscriber or it belongs to the rep-set R_d). Consequently, the potential neighbor ratio is always the biggest for lightweight nodes selected as rep-nodes. For such nodes, $pn\text{-ratio}(v) = \frac{|L|}{|V| \cdot p} + \frac{|B|}{|V|} + \frac{|R|}{|V|} \cdot \frac{p-1}{p}$.

We extend the definition of a potential neighbor ratio to apply to the entire node set:

$$pn\text{-ratio}(V) = \max_{v \in V} pn\text{-ratio}(v) = \frac{|L|}{|V| \cdot p} + \frac{|B|}{|V|} + \frac{|R|}{|V|} \cdot \frac{p-1}{p} \quad (9)$$

Equation 9 shows that DCBR-M has improved pn -ratio compared to any centralized algorithm (such as GM-M). This is further confirmed by our experiments in Section VI.

V. SELECTING PARAMETERS

This section discusses how to choose optimal parameter values for our DCBR-M algorithm. Algorithm 7 is parametrized with (1) the bulk subscriber threshold η , (2) the coverage factor λ , and (3) the number of partitions for lightweight nodes p . The choice of values for these parameters substantially affects the algorithm’s behavior. It is therefore essential to identify a combination of them that leads to satisfactory performance. First, we pick reasonable values for η and λ for typical pub/sub workloads; then, we provide a numerical method to determine the optimal value of p .

The selection of the bulk subscriber threshold η exhibits the tradeoff between maximal overlay degree and running time: Small threshold values cause all nodes to be treated as bulk subscribers thereby favoring the degree over running time and making the overall performance very similar to that of GM-M. On the other hand, large threshold values favor the running time and pn -ratio at the expense of increased overlay degree. Fortunately, even relatively small threshold values result in small bulk subscriber sets for typical pub/sub workloads that follow the “Pareto 80–20” rule, as discussed in Section IV-B. In our implementation, we sort the subscribers by subscription size and choose η that causes $\leq 20\%$ of the nodes to be considered bulk subscribers.

The coverage factor selection exhibits a similar tradeoff: If we choose the coverage factor to be as large as the size $|L|/p$ of the partitions, then the behavior of Algorithm 7 becomes identical to that of Algorithm 5. On the other hand, $\lambda = 1$ minimizes the size of rep-sets, pn -ratio, and running time but leads to a severe impact on the node degree. According to our experiments in Section VI, an increase in λ beyond 3 only marginally improves the node degree, even for large partitions. The rep-sets for $\lambda = 3$ are significantly smaller for such large partitions than partitions themselves so that we choose 3 as the default value for λ .

The tradeoff in the selection of the number p of partitions is more complex. When p is as large as $|L|$, the performance is dominated by the invocations of the `combineB&LReps()` function in the *combine* phase. As p decreases, the effect of executing the GM-M algorithm at the *conquer* phase becomes more and more pronounced, both with respect to the degree and the running time. When we use a very small number of partitions, it starts to dominate the running time assuming $|B|$ is relatively small. Therefore, we need to find intermediate values of p that minimize the running time and pn -ratio.

The bound on the running time \mathbb{T}_{DCBR} is established by Lemma 12 and Equation 9 for centralized and distributed implementations, respectively. Since the bound on \mathbb{T}_{DCBR} depends on $|R|$, which is difficult to assess analytically, we use an adaptive way for selecting p . Since partitioning the nodes and computing the rep-sets is relatively cheap, we try partitioning for different p values. Each time, we only

TABLE I
ALGORITHMS FOR SOLVING THE MINMAX-TCO PROBLEM

GM-M	Greedy Merge algorithm for MinMax
DCB-M	Divide-and-Conquer with Bulk and Lightweight Nodes
DCBR-M	Divide-and-Conquer with Bulk and Lightweight Rep-nodes
RingPT	Ring-per-topic algorithm
TCO_{ALG}^*	The TCO produced by ALG
D_{ALG}	Maximum node degree in TCO_{ALG}
d_{ALG}	Average node degree in TCO_{ALG}
\mathbb{T}_{ALG}	Running time of ALG

* ALG is GMM, DCB, DCBR or RingPT for the GM-M, DCB-M, DCBR-M or RingPT algorithms, respectively.

compute the rep-sets and thus obtain $|R|$ without running the expensive calculation of inner and outer edges. Then, we use fast numerical methods to approximately determine the value of p that minimizes \mathbb{T}_{DCBR} . We also apply the same technique to Equation 9 in order to determine the value of p that is optimal for pn -ratio(V).

VI. EVALUATION

We implemented all algorithms described in this paper in Java and compared them under various experimental conditions. Table I summarizes the algorithms evaluated. We use the GM-M algorithm as a baseline because it produces the lowest maximum node degree of all known algorithms that run in polynomial time. To be precise, we are using a faster implementation of the GM-M algorithm described in [26] both for baseline GM-M and as a building block for DCB-M and DCBR-M. This faster implementation produces exactly the same overlay as the original one in [8] at a lower runtime cost by manipulating data structures more efficiently. Our divide-and-conquer design is orthogonal to the data structures used in the algorithm: Our faster version still has prohibitively high running time without divide-and-conquer. In fact, the speedup of DCBR-M compared to GM-M would have been even more significant for the slower, original GM-M implementation. However, using a faster implementation allows us to run comparative experiments on a larger scale.

In the experiments, we use the following ranges for the input instances: $|V| \in [1\,000, 8\,000]$ and $|T| \in [100, 1\,000]$. We define the average node subscription size, minimum subscription size, and the maximum subscription size as follows: $\overline{|T_v|} = \frac{\sum_{v \in V} |T_v|}{|V|}$, $|T_v|_{\min} = \min_{v \in V} \{|T_v|\}$, $|T_v|_{\max} = \max_{v \in V} \{|T_v|\}$. We used $|V| = 4\,000$, $|T| = 200$, and $\overline{|T_v|} = 50$ (with $|T_v|_{\min} = 10$, $|T_v|_{\max} = 90$) to generate the input workloads for most of the experiments unless specified otherwise. Each topic $t_i \in T$ is associated with probability q_i , $\sum_i q_i = 1$, so that each node subscribes to t_i with a probability q_i . The value of q_i is distributed according to either a uniform, a Zipf (with $\alpha = 2.0$), or an exponential distribution. According to [21], these distributions are representative of actual workloads used in industrial pub/sub systems today. The Zipf distribution is chosen because [3] shows it faithfully describes the feed popularity distribution in RSS. The exponential distribution is used by stock-market monitoring engines in [28] for the study of stock popularity in the New York Stock Exchange (NYSE). The η , p , and λ parameters are selected as described in Section V.

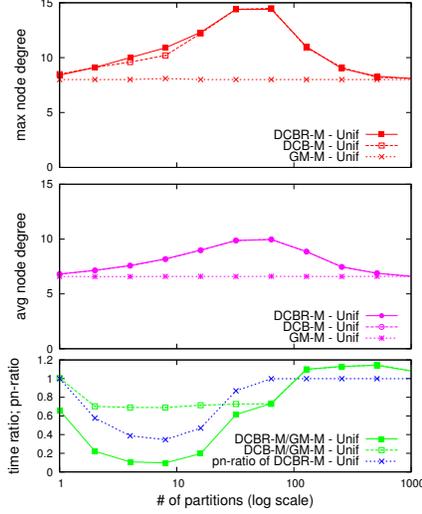


Fig. 4. DCBR-M parameterized with different p

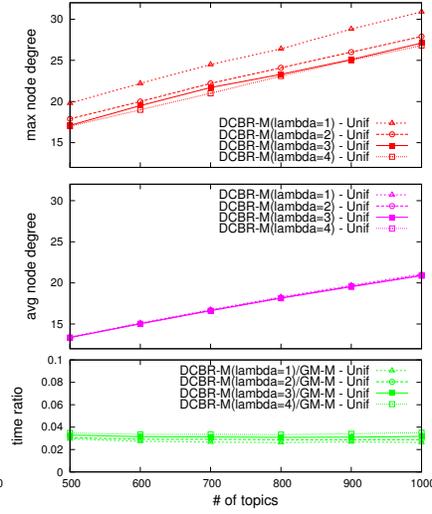


Fig. 5. DCBR-M parameterized with different λ

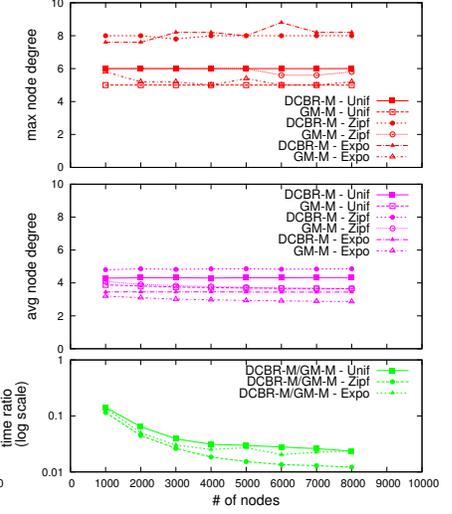


Fig. 6. DCBR-M under different distributions

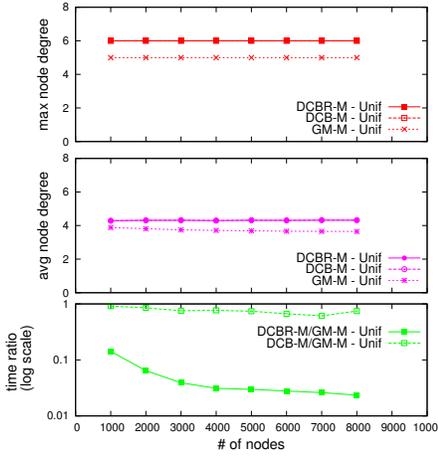


Fig. 7. DCBR-M as $|V|$ increases

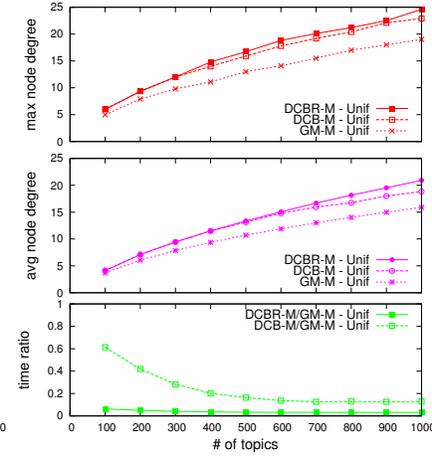


Fig. 8. DCBR-M as $|T|$ increases

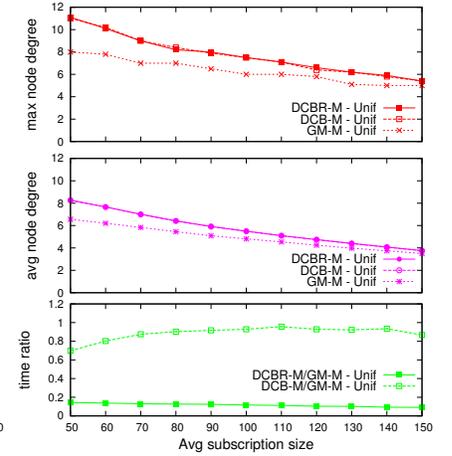


Fig. 9. DCBR-M as $|\overline{T}_v|$ increases

A. Partitioning for Lightweight Nodes

TABLE II
RANDOM PARTITIONING UNDER EXPONENTIAL DISTRIBUTION

	mean	min	max	variance
D_{DCBR}	10.798	9	13	0.468
D_{DCB}	10.793	9	13	0.461
D_{GMM}	8.4425	7	12	0.348
d_{DCBR}	4.499	4.394	4.586	0.00122
d_{DCB}	4.499	4.392	4.59	0.00124
d_{GMM}	3.93	3.86	4.02	0.000747
$\frac{T_{DCBR}}{T_{GMM}}$	0.130	0.111	0.161	0.000678
$\frac{T_{DCB}}{T_{GMM}}$	0.833	0.787	0.883	0.000225
$ R / L $	0.0748	0.0438	0.0977	0.0000891

a) *Random partitioning for lightweight nodes:* We first evaluate the effects of random partitioning of lightweight nodes for the DCB-M and DCBR-M algorithms. We run the algorithm 400 times for the same settings (namely, the default experimental settings discussed above except $|V| = 1000$, under three different distributions) so that the only difference between different runs is due to the random node interest generation according to the given distribution parameters and

due to random node partitioning. The statistics pertaining to maximum node degree, average node degree, running time ratio, and the ratio of nodes selected as rep-nodes are reported in Table II. As the table illustrates, under the exponential distribution, all the values are quite *stable* with negligible variance across different experiments; and the variance is even more insignificant under less skewed distributions. Besides, the results validate our assumption that $|R| \ll |L|$ (with $\lambda = 3$). We conclude that when the number p of partitions is *reasonable*, random partitioning of lightweight nodes is an efficient and robust way to implement the *divide* phase of DCBR-M. Furthermore, it results in small rep-sets, which is vital to the performance of the DCBR-M algorithm.

b) *Impact of p :* Given an instance $I(V, T, Int)$ where $|V| = 1000$, DCBR-M and DCB-M are executed with all possible p values ranging from 1 to $|L|$. Fig. 4 shows that under different values of p , TCO_{DCBR} and TCO_{DCB} have similar maximum and average node degrees, which are slightly higher than those of TCO_{GMM} . However, DCBR-M runs substantially faster than DCB-M when p is set appropriately.

We already know that DCBR-M and DCB-M exhibit identical behavior to GM-M when $p = |L|$. Fig. 4 shows that as the number of partitions p grows from 1 to $|L|$, $D_{\text{DCBR}} (\approx D_{\text{DCB}})$ first *increases* gradually and then it starts to *decrease* until it becomes equal to $D_{\text{GM-M}}$. Note that the D_{DCBR} never moves far from the horizon line of D_{GMM} : It is always smaller than 6.5. It is less than 2.8 for the value of p that minimizes \mathbb{T}_{DCBR} .

While \mathbb{T}_{DCB} stays close to \mathbb{T}_{GMM} for all values of p , \mathbb{T}_{DCBR} first slides *down* sharply and then climbs *up* as p increases. The range of p values for which this phenomenon occurs is small and relatively close to 0. (It touches the lowest point when p is around 10.) Furthermore, the $\mathbb{T}_{\text{DCBR}}/\mathbb{T}_{\text{GMM}}$ ratio follows the same trend as *pn-ratio* for DCBR-M, which is compatible with our discussion of choosing p in Section V.

B. Impact of Coverage Factor

Here, we explored the impact of λ on the output and performance of the DCBR-M algorithm. Given an input $I(V, T, Int)$ where $|V| = 2000$ and $|T| \in [500, 1000]$, the DCBR-M algorithm is evaluated for four different values of the coverage factor ($\lambda = 1, 2, 3, 4$). As Fig. 5 shows, under uniform distribution, as λ *increases*, D_{DCBR} *decreases*. The differences in maximum node degrees also *decrease* with successive coverage factors, i.e., $D_{\text{DCBR}}|_{\lambda=k-1} - D_{\text{DCBR}}|_{\lambda=k} > D_{\text{DCBR}}|_{\lambda=k} - D_{\text{DCBR}}|_{\lambda=k+1}$. More specifically, compared to $D_{\text{DCBR}}|_{\lambda=2} - D_{\text{DCBR}}|_{\lambda=3} \approx 0.71$ on average, $D_{\text{DCBR}}|_{\lambda=1} - D_{\text{DCBR}}|_{\lambda=2}$ is noticeable, which could be as much as ≥ 3 for most cases. When $\lambda \geq 3$, the difference is insignificant (≤ 0.32 on average). Meanwhile, under all coverage factors that were tested, DCBR-M runs remarkably faster as compared to GM-M ($\mathbb{T}_{\text{DCBR}} \leq 3\% \mathbb{T}_{\text{GMM}}$ on average). Also, \mathbb{T}_{DCBR} slightly *increases* as λ *increases*, because *pn-ratio* *increases* as a result. However, differences among the running time cost for different coverage factors are insignificant: $\frac{\mathbb{T}_{\text{DCBR}}|_{\lambda=4}}{\mathbb{T}_{\text{GMM}}} - \frac{\mathbb{T}_{\text{DCBR}}|_{\lambda=1}}{\mathbb{T}_{\text{GMM}}} \leq 0.69\%$ on average.

This experiment confirms the validity of choosing a relatively small integer as coverage factor in Section V, both in terms of node degree and running time cost.

C. Effects Under Different Distributions

We now consider DCBR-M's behavior under different input instances. We first provide an overview of the overall DCBR-M performance under different typical distributions, and then we analyze how DCBR-M is affected by various aspects of the input.

Fig. 6 depicts that under different distributions, DCBR-M produces high-quality TCOs in terms of maximum and average node degrees, which are slightly higher than D_{GMM} and d_{GMM} , respectively. However, the differences are insignificant: $D_{\text{DCBR}} - D_{\text{GMM}} \leq 2.0$, $d_{\text{DCBR}} - d_{\text{GMM}} \leq 0.70$ on average.

Although DCBR-M and GM-M produce quite close maximum and average node degrees, DCBR-M runs considerably faster than GM-M: $\mathbb{T}_{\text{DCBR}} \leq 4.0\% \cdot \mathbb{T}_{\text{GMM}}$ on average. As the number of nodes *increases*, D_{DCBR} and d_{GMM} remain steadily *low* while the running time ratio $\mathbb{T}_{\text{DCBR}}/\mathbb{T}_{\text{GMM}}$ *decreases*

considerably. This further attests to DCBR-M's scalability with respect to the number of nodes in the network.

The maximum node degrees tend to be a bit more *fluctuating* under Zipf and exponential distributions compared to those under the uniform distribution. This could be explained by the slightly higher variance under skewed distributions, as presented in Table II. Although skewed distributions are more sensitive to the variations in the input, the maximum and average node degrees always stay low, even in the worst cases.

D. Impact of the number of nodes

We now demonstrate DCBR-M's scalability with respect to different input parameters. In the rest of the section, while we report on results and analysis for all distributions in the text, we only show figures for the uniform topic popularity distribution due to space limit. See [26] for additional results.

Fig 7 depicts the comparison between DCBR-M, DCB-M and GM-M as the number of nodes increases where $|T| = 100$. The figure shows that DCBR-M and DCB-M output similar TCOs with regard to maximum and average node degrees, but DCBR-M runs considerably faster. Under the uniform distribution, for example, \mathbb{T}_{DCBR} is on average 4.79% of \mathbb{T}_{GMM} while \mathbb{T}_{DCB} is as much as 75.7% of \mathbb{T}_{GMM} . Additionally, DCBR-M gains more speedup with the increase in the number of nodes compared to the other algorithm.

E. Impact of the number of topics

Fig. 8 depicts how DCBR-M and DCB-M perform compared to GM-M when we vary the number of topics. As the figure shows, under the uniform distribution, the maximum and average node degrees of all three algorithms *increase* for a higher number of topics. This is because *increasing* the number of topics leads to *reduced* correlation among subscriptions. However, the increase is slow paced and the difference $D_{\text{DCBR}} - D_{\text{GMM}}$ remains insignificant: 3.51 for the uniform, 4.46 for the Zipf, and 3.46 for the exponential distribution on average.

The running time ratio of DCBR-M to GM-M slightly *increases* as the number of topics *increases*, yet this effect is insignificant: \mathbb{T}_{DCBR} is less than 3.8% of \mathbb{T}_{GMM} on average.

F. Impact of the average subscription size

Fig. 9 depicts how the node subscription size affects the DCBR-M and DCB-M algorithms. We set $|V| = 1000$, $|T| = 400$, and $\overline{|T_v|}$ varies from 50 to 150.

The figure shows that under the uniform distribution, DCBR-M and DCB-M produces quite close TCOs in terms of both maximum and average node degrees. As the subscription size *increases*, D_{DCBR} and D_{GMM} *decrease*, and the difference of $(D_{\text{DCBR}} - D_{\text{GMM}})$ *shrinks*. d_{DCBR} follows the same trend.

This *decrease* occurs because the growth of $\overline{|T_v|}$ causes *increased* correlation across the subscriptions. Upon bigger correlation, an edge addition to the overlay reduces a higher number of topic-connected components on average because the nodes share more comment interests. Therefore, a smaller number of edge additions are required before the overlay becomes topic-connected.

The ratio of \mathbb{T}_{DCBR} to \mathbb{T}_{GMM} also *decreases* with the *increase* of $|T_v|$. In both algorithms, an edge addition causes a higher number of updates to topic-connected components for bigger $|T_v|$ (Lines 8–9 in Algorithm 2). Yet, this effect has less influence on \mathbb{T}_{DCBR} compared to \mathbb{T}_{GMM} since each update in DCBR-M affects a smaller portion of edges. Unlike \mathbb{T}_{DCBR} , however, \mathbb{T}_{DCB} does not gain significant speedup.

G. Comparison with Ring-Per-Topic

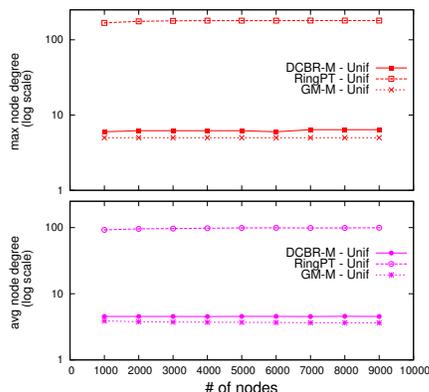


Fig. 10. DCBR-M vs. GM-M vs. RingPT

Finally, we compare the maximum and average node degrees produced by DCBR-M GM-M and RingPT. RingPT is an algorithm that mimics the common practice of building a separate overlay for each topic (usually a tree but we use a ring that has the same average node degree). According to RingPT, all the nodes interested in the same topic form a ring for that topic, after which rings for different topics are merged into a single overlay. $|T|$ is set to 100 in this experiment. As Fig. 10 shows, D_{DCBR} and D_{GMM} are quite close (the average difference is 1.22), but the maximum node degree of RingPT exceeds D_{DCBR} by a factor of approximately 30. This demonstrates the general significance of overlay construction algorithms for pub/sub.

VII. CONCLUSIONS

This paper focuses on a number of design objectives for the MinMax-TCO problem that are central to creating a practical solution. We have designed the DCBR-M algorithm which is capable of constructing a low fan-out TCO, while being significantly more efficient than previously known solutions. Numerical techniques can be employed to effectively obtain a good combination of parameters which adapts to various inputs and guarantees the output and the performance of the algorithm. The algorithm is thoroughly examined via a comprehensive experimental analysis, which demonstrates the scalability of DCBR-M under different distributions as the number of nodes, the number of topics, and the subscription size increase.

REFERENCES

[1] J. Reumann, “Pub/Sub at Google,” lecture & Personal Communications at EuroSys & CANOE Summer School, Oslo, Norway, Aug’09.

[2] “Tibco rendezvous,” <http://www.tibco.com>.

[3] H. Liu, V. Ramasubramanian, and E. G. Sifer, “Client behavior and feed characteristics of RSS, a publish-subscribe system for web micronews,” in *IMC’05*.

[4] M. Petrovic, H. Liu, and H.-A. Jacobsen, “G-ToPSS: fast filtering of graph-based metadata,” in *WWW’05*, 2005.

[5] G. Li, V. Muthusamy, and H.-A. Jacobsen, “A distributed service oriented architecture for business process execution,” *ACM TWEB*, 2010.

[6] B. F. Cooper, R. Ramakrishnan, U. Srivastava, A. Silberstein, P. Bohannon, H.-A. Jacobsen, N. Puz, D. Weaver, and R. Yerneni, “Pnuts: Yahoo!’s hosted data serving platform,” *Proc. VLDB Endow.*, 2008.

[7] G. Chockler, R. Melamed, Y. Tock, and R. Vitenberg, “Constructing scalable overlays for pub-sub with many topics: Problems, algorithms, and evaluation,” in *PODC’07*.

[8] M. Onus and A. W. Richa, “Minimum maximum degree publish-subscribe overlay network design,” in *INFOCOM’09*.

[9] —, “Parameterized maximum and average degree approximation in topic-based publish-subscribe overlay network design,” in *ICDCS’10*.

[10] M. A. Jaeger, H. Parzyjegl, G. Mühl, and K. Herrmann, “Self-organizing broker topologies for publish/subscribe systems,” in *SAC’07*.

[11] R. Baldoni, R. Beraldi, L. Querzoni, and A. Virgillito, “Efficient publish/subscribe through a self-organizing broker overlay and its application to SIENA,” *Comput. J.*, vol. 50, no. 4, 2007.

[12] C. Chen, H.-A. Jacobsen, and R. Vitenberg, “Divide and conquer algorithms for publish/subscribe overlay design,” in *ICDCS’10*.

[13] D. Liben-Nowell, H. Balakrishnan, and D. Karger, “Analysis of the evolution of peer-to-peer systems,” in *PODC*, 2002.

[14] E. De Santis, F. Grandoni, and A. Panconesi, “Fast low degree connectivity of ad-hoc networks via percolation,” in *ESA’07*.

[15] L. C. Lau, J. S. Naor, M. R. Salavatipour, and M. Singh, “Survivable network design with degree or order constraints,” in *Proc. ACM STOC’07*.

[16] D. Tam, R. Azimi, and H.-A. Jacobsen, “Building content-based publish/subscribe systems with distributed hash tables,” in *DBISP2P’03*.

[17] M. Castro, P. Druschel, A.-M. Kermarrec, and A. Rowstron, “SCRIBE: A large-scale and decentralized application-level multicast infrastructure,” *JSAC*, 2002.

[18] G. Li, V. Muthusamy, and H.-A. Jacobsen, “Adaptive content-based routing in general overlay topologies,” in *Middleware’08*.

[19] F. Araujo, L. Rodrigues, and N. Carvalho, “Scalable QoS-based event routing in publish-subscribe systems,” in *NCA’05*.

[20] S. Girdzijauskas, G. Chockler, Y. Vigfusson, Y. Tock, and R. Melamed, “Magnet: practical subscription clustering for internet-scale publish/subscribe,” in *DEBS’10*.

[21] G. Chockler, R. Melamed, Y. Tock, and R. Vitenberg, “Spidercast: A scalable interest-aware overlay for topic-based pub/sub communication,” in *DEBS’07*.

[22] R. Baldoni, R. Beraldi, V. Quema, L. Querzoni, and S. Tucci-Piergiovanni, “TERA: topic-based event routing for peer-to-peer architectures,” in *DEBS’07*.

[23] E. Baehni, P. Eugster, and R. Guerraoui, “Data-aware multicast,” in *DSN’04*.

[24] R. Chand and P. Felber, “Semantic peer-to-peer overlays for publish/subscribe networks,” in *EUROPAR’05*.

[25] S. Voulgaris, E. Rivire, A.-M. Kermarrec, and M. V. Steen, “Sub-2-Sub: Self-organizing content-based publish subscribe for dynamic large scale collaborative networks,” in *IPTPS’06*.

[26] C. Chen, R. Vitenberg, and H.-A. Jacobsen, “Scaling construction of low fan-out overlays for topic-based publish/subscribe systems,” U. of Toronto & U. of Oslo, Tech. Rep., 2010, <http://msrg.org/papers/TRCJV-DCBRM-2010>.

[27] D. Peleg, G. Schechtman, and A. Wool, “Approximating bounded 0-1 integer linear programs,” in *Theory of Computing and Systems*, 1993.

[28] Y. Tock, N. Naaman, A. Harpaz, and G. Gershinsky, “Hierarchical clustering of message flows in a multicast data dissemination system,” in *IASTED PDCS*, 2005.