

Relevance Matters: Capitalizing on Less

Top-k Matching in Publish/Subscribe

Mohammad Sadoghi¹, Hans-Arno Jacobsen²

Department of Computer Science, University of Toronto

¹mo@cs.toronto.edu ²jacobsen@eecg.toronto.edu

Abstract—The efficient processing of large collections of Boolean expressions plays a central role in major data intensive applications ranging from user-centric processing and personalization to real-time data analysis. Emerging applications such as computational advertising and selective information dissemination demand determining and presenting to an end-user only the most relevant content that is both user-consumable and suitable for limited screen real estate of target devices. To retrieve the most relevant content, we present BE*-Tree, a novel indexing data structure designed for effective hierarchical top-k pattern matching, which as its by-product also reduces the operational cost of processing millions of patterns.

To further reduce processing cost, BE*-Tree employs an adaptive and non-rigid space-cutting technique designed to efficiently index Boolean expressions over a high-dimensional continuous space. At the core of BE*-Tree lie two innovative ideas: (1) a bi-directional tree expansion build as a top-down (data and space clustering) and a bottom-up growths (space clustering), which together enable indexing only non-empty continuous sub-spaces, and (2) an overlap-free splitting strategy. Finally, the performance of BE*-Tree is proven through a comprehensive experimental comparison against state-of-the-art index structures for matching Boolean expressions.

I. INTRODUCTION

The efficient processing of Boolean expressions is a key functionality required by major data intensive applications ranging from user-centric processing and personalization to real-time data analysis. Among user-centric applications, there are targeted Web advertising, online job sites, and location-based services; common to all are specifications (e.g., advertising campaigns, job profiles, service descriptions) modeled as Boolean expressions (subscriptions) and incoming user information (e.g., user profiles) modeled as events [1], [2], [3]. Among real-time data analysis applications, there are (complex) event processing, intrusion detection, and computational finance; again, common among them are predefined set of patterns (e.g., attack specifications and investment strategies) modeled as Boolean expressions (subscriptions) and streams of incoming data (e.g., stock feeds) modeled as events [4], [5], [6], [2].

Unique to user-centric processing and personalization are strict requirements to determine only the most relevant content that is both user-consumable and suitable for the often limited screen real estate of client devices. Shared among user-centric processing and real-time analysis applications are two requirements: (1) scaling to millions of patterns and specifications (expressions) and (2) matching latency constraints in the sub-second range.

To put into context the importance of finding the most relevant Boolean expression, we extract real scenarios from

prominent user-centric applications such as targeted Web advertising (e.g., Google, Microsoft, Yahoo!) and online job sites (e.g., Monster).

In Web advertising, a demographic targeting service specifies constraints such as [*credit-score* > 650 (*wt* = 0.4), *num-visits* > 10 (*wt* = 0.2), *age* BETWEEN {22, 36} (*wt* = 0.4)] while an incoming user's profile includes information such as [*credit-score* = 732 (*wt* = 0.2), *num-visits* = 17 (*wt* = 0.2), *age* = 27 (*wt* = 0.6)]. Thus, only the most relevant ads for the user are displayed; the ad relevance is computed based on a scoring function that takes as inputs the constraints' weights (*wt*) assigned by user and advertiser. Similarly, in online job sites, an employer submits the job details, [*category* = 'green jobs' (*wt* = 0.2), *hours/week* > 15 (*wt* = 0.5), *hourly-rate* < 45 (*wt* = 0.3)], and when a job seeker visits with the profile, [*category* = 'green jobs' (*wt* = 0.4), *hours/week* = 20 (*wt* = 0.1), and *hourly-rate* = 30 (*wt* = 0.5)], she is connected to the most relevant employer [1], [2], [3]; again the relevance is computed based on the weights of the constraints. These scenarios demand the indexing of millions of Boolean expressions while having matching latency requirements in the sub-second range to avoid the deterioration of user-perceived latency when browsing and shopping online.

In these examples, advertisers' and employers' needs and users' profiles are captured using multiple attributes (i.e., dimensions in space) such as *age* and *credit-score*, in which the importance of each attribute is explicitly captured using the attributes' weight. Additionally, each attribute is constrained to a range of values on either continuous or discrete domains using an operator. The resulting quadruple formed by an attribute, an operator, a range of values, and a predicate weight is referred to as a *Boolean predicate*. A conjunction of Boolean predicates, here representing both the advertiser's and employer's needs and the user's profiles, constitutes a *Boolean expression (BE)*.

In short, we extract four requirements of paramount importance from applications that rely on processing Boolean expressions. First and foremost, the index must enable top-k matching to quickly retrieve only the most relevant expressions, which is not addressed by prior Boolean expression matching approaches [4], [5], [7]. This relevance computation must be based on a generic scoring function, which is also not addressed by prior approaches that model relevance identification as part of matching [1]. Second, the index must support predicates with an expressive set of operators over a continuous domain, which, again, is not supported by prior matching approaches [5], [2], [7]. Third, the index must enable

dynamic insertion and deletion of expressions, often disregarded as requirement in matching algorithm designs [4], [2]. Finally, the index must employ a dynamic structure that adapts to changing workload patterns and expression schemata, also often disregarded in matching algorithm designs [4], [2].

We designed BE*-Tree to specifically tackle these challenges. In particular, BE*-Tree is geared towards determining the most relevant patterns, in which top-k processing is treated as first class citizen. BE*-Tree introduces a hierarchical top-k processing scheme that differs from existing work that assumes a static and a flat structure [2]. Additionally, BE*-Tree achieves scalability by overcoming the curse of dimensionality through a non-rigid space-cutting technique while attaining expressiveness and dynamism without restricting the space to only a discrete and finite domain as prevalent in [2], [7]. Notably, with respect to scalability, we solve the two critical properties common to most high-dimensional indexing: (1) avoiding indexing non-empty space (cf. R-tree [8]) (2) minimizing overlap (cf. R*-tree [9]) and coverage [8]. These properties are tackled by proposing a bi-directional tree expansion: a top-down (data and space clustering) and a bottom-up (space clustering) growths process, which together enable indexing only non-empty continuous sub-spaces and adaptation to workload changes; and a splitting strategy to systematically produce and maintain overlap-free subspaces for holding expressions.

Our contributions in this paper are four-fold:

- 1) We generalize the BE matching problem to subspace matching (Section III),
- 2) We introduce BE*-Tree tailored to top-k processing for efficiently determining the most relevant matches among large numbers of potential matches of millions of subscriptions defined over thousands of dimensions (Sections IV, VI),
- 3) We develop a novel self-adjusting mechanism using bi-directional tree expansion that continuously adapts as the workload changes, and propose a novel overlap-free splitting strategy (Section V),
- 4) Finally, we present a comprehensive evaluation framework that benchmarks state-of-the-art matching algorithms, specifically emphasizing top-k processing: SCAN [10], BE-Tree [7], and k -index [2] (Section VII).

II. RELATED WORK

Generally speaking the problems related to indexing Boolean expressions have been studied in different contexts including database (multi-dimensional indexing [8], [9], [11]) and top-k querying [12], [13]) and publish/subscribe (matching techniques [10], [4], [5], [6], [2], [3], [7]).

In general, the database literature focuses on much lower dimensionality, while we target a dimensionality in order of thousands, which is orders of magnitude larger than capabilities of existing database indexing structures [11]. The database top-k processing [12], [13] differs with our proposed top-k model in an important respect: our top-k model solves the reverse problem. In the database context, top-k querying means finding the most relevant tuples (events) for a given query (subscription). But in our context, top-k matching means

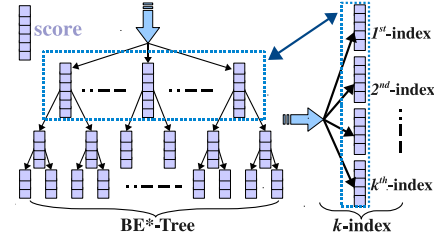


Fig. 1. Top-k Structural Comparisons

finding the most relevant subscriptions (queries) for a given event (tuple).

In publish/subscribe, two main categories of matching algorithms have been proposed: counting-based [10], [5], [2] and tree-based [4], [6], [7] approaches. These approaches can further be classified as either key-based, in which for each expression a set of predicates are chosen as identifier [5], [7], or as non-key based [10], [6], [2]. Alternatively, these approaches can also be classified either as rigid- [10], [5], [2], [7] and non-rigid-based (BE*-Tree) approaches. What constitutes rigidity is a pre-determined clustering of the space. For example, for a skewed distribution in which all data is clustered in a small subspace, rigid clustering continues splitting the space in half, irrespective of where the data is actually located, and is unable to avoid indexing the empty subspace.

A closer look manifests that counting-based methods primarily aim to minimize the number of predicate evaluations by constructing an inverted index over all unique predicates resulting in a rigid clustering. The two most efficient counting-based algorithms are Propagation [5], a key-based method, and the k -index [2], a non-key-based method. However, BE*-Tree, unlike Propagation, supports promoting non-equality predicates as key and supports continuous domains. Furthermore, BE*-Tree unlike k -index, is a dynamic structure that supports both discrete and continuous domains.

Likewise, tree-based methods are primarily designed to reduce predicate evaluations and to recursively divide search space by eliminating expressions on encountering unsatisfiable predicates. The first major tree-based approach, Gryphon, is a static, a non-key based method [4], which is shown to be effective only for equality predicates [4]. The latest tree-based structure is BE-Tree, a key-based approach, that introduces a two-phase space-cutting abstraction for supporting workload changes and overcoming the curse of dimensionality [7]. BE-Tree leverages an effective rigid clustering that is inherently designed for finite and discrete domains.

In comparison, BE*-Tree develops a non-rigid clustering through a self-adjusting mechanism using the bi-directional tree expansion and the overlap-free splitting strategy that together let BE*-Tree adapt to skewed workloads, cluster only non-empty subspaces, and improve search space pruning. Above all, BE*-Tree enables the top-k relevance-match-determination model. Although BE*-Tree addresses the rigidity and top-k limitations of BE-Tree, yet BE-Tree has two notable advantages over BE*-Tree: it is algorithmically easier to implement and robust w.r.t. insertion sequence (a desirable property also lacking in R-Tree family of indexes [11]).

A fresh look at enhancing existing pub/sub matching algorithms is to leverage top-k processing techniques, which not only concentrate on determining the most relevant matches, but can also improve matching performance. An early top-k model is presented in [1]; however, this model is based on a fixed and predetermined scoring function, i.e., the score for each expression is computed independent of the incoming event. In addition, this approach is an extension of the *R*-Tree, the interval tree, or the segment tree structure; as a result, it has difficulties scaling beyond a few dimensions [1]. In contrast, a scalable top-k model, but based on a static and flat structure, with a generic scoring function, which also takes the event into consideration, is introduced in *k*-index [2]. BE*-Tree's top-k model not only scales to thousands of dimensions but introduces a hierarchical top-k model. Our top-k model supports a generic scoring function, that achieves an effective pruning power for determining the most relevant matches.

BE*-Tree's hierarchical top-k structure differs from the flat structure of *k*-index [2]. This key structural differentiation is depicted in Figure 1. A hierarchical model achieves enhanced pruning because it refines upper bound scores as more expressions are seen (as traversing down the tree). Such pruning is not possible with a flat structure, in which the upper bound score is computed once at the top-layer because as more expressions are seen, the new information is not reflected back in the original top-layer scores.

III. LANGUAGE AND DATA MODEL FORMALISM

In this section, we define our Boolean expression language and spatial event data model followed by our (top-k) matching semantics.

A. Notation

Given an n -dimensional space \mathbb{R}^n , we define the projection of \mathbb{R}^n onto \mathbb{R}^k as a k -dimensional subspace, denoted by $\pi_{d_1 \dots d_k}(\mathbb{R}^n) = \mathbb{R}^k$, where $\pi_{d_1 \dots d_k} : \mathbb{R}^n \rightarrow \mathbb{R}^k$, $k \leq n$, and each $d_i \in \{d_1 \dots d_k\}$ represents the i^{th} dimension in \mathbb{R}^k and corresponds to the j^{th} dimension in \mathbb{R}^n ; for ease of notation, we define the identity projection as $\pi_{\mathcal{I}}(\mathbb{R}^n) = \mathbb{R}^n$. In addition, we define a k -dimensional bounding box \mathcal{B}^k over \mathbb{R}^k as

$$\mathcal{B}^k = [\min_1, \max_1] \times \dots \times [\min_k, \max_k].$$

Let $\xi_j(\mathcal{B}^k) = [\min_i, \max_i]$ be the i^{th} boundary in \mathcal{B}^k defined over the i^{th} dimension in \mathbb{R}^k which corresponds to the j^{th} dimension in \mathbb{R}^n . Let $\chi_i(\mathcal{B}^k)$ be the center of the i^{th} boundary in \mathcal{B}^k . Furthermore, let $\lambda_i(\mathcal{B}^k)$ be the length of the i^{th} boundary in \mathcal{B}^k given by

$$\lambda_i(\mathcal{B}^k) = \max_i - \min_i.$$

Lastly, let $\mu_i(\mathcal{B}^k) = \min_i$ and $M_i(\mathcal{B}^k) = \max_i$ represent the minimum and the maximum value of the i^{th} boundary of \mathcal{B}^k , respectively.

B. Expression Language and Subspace Model

We support a rich Boolean expression language that unifies the subscription language and the event data model. This generalization gives rise to more expressive matching semantics while still encompassing the traditional pub/sub matching problem.

TABLE I
PREDICATE MAPPING

$P^{\text{attr, opt, val, wt}}(x)$	$\gamma(P^{\text{attr, opt, val, wt}}(x))$
$\text{attr} < v_1$	$[-\infty, v_1 - \epsilon]$
$\text{attr} \leq v_1$	$[-\infty, v_1]$
$\text{attr} = v_1$	$[v_1, v_1]$
$\text{attr} \neq v_1$	$[-\infty, \infty]$
$\text{attr} \geq v_1$	$[v_1, \infty]$
$\text{attr} > v_1$	$[v_1 + \epsilon, \infty]$
$\text{attr} \in \{v_1, \dots, v_m\}$	$[v_1, v_m]$
$\text{attr} \notin \{v_1, \dots, v_m\}$	$[-\infty, \infty]$
$\text{attr BETWEEN } v_1, v_2$	$[v_1, v_2]$

In our model, a Boolean expression is a conjunction of Boolean predicates. A predicate is a quadruple: an attribute that uniquely represents a dimension in \mathbb{R}^n ; an operator (e.g., relational operators ($<, \leq, =, \neq, \geq, >$), set operators (\in, \notin), and the SQL BETWEEN operator); a set of values (for discrete domains) or a range of values (for continuous domains); and an assigned predicate weight, denoted by $P^{\text{attr, opt, val, wt}}(x)$ or more concisely as $P(x)$. A predicate either accepts or rejects an input x such that $P^{\text{attr, opt, val, wt}}(x) : x \rightarrow \{\text{True}, \text{False}\}$, where $x \in \text{Dom}(P^{\text{attr}})$ and P^{attr} is the predicate's attribute. Formally, a Boolean expression Ω is defined over \mathbb{R}^n as follows:

$$\Omega = \{P_1^{\text{attr, opt, val, wt}}(x) \wedge \dots \wedge P_k^{\text{attr, opt, val, wt}}(x)\},$$

where $k \leq n$; $i, j \leq k$, $P_i^{\text{attr}} = P_j^{\text{attr}}$ iff $i = j$.

We extended the semantics of projection to a Boolean expression in order to enable projecting out predicates associated with certain dimensions

$$\pi_{d_1 \dots d_h}(\Omega) = \{\Omega' \mid \Omega' = P_1^{\text{attr, opt, val, wt}}(x) \wedge \dots \wedge P_h^{\text{attr, opt, val, wt}}(x),$$

$$\forall d_i \in \{d_1 \dots d_h\}, P_i^{\text{attr}} = d_i, P_i(x) \in \Omega\}.$$

Now, we are in a position to formalize the predicate mapping¹ and, ultimately, to (approximately) represent an expression as a k -dimensional bounding box.

A predicate $P^{\text{attr, opt, val, wt}}(x)$ is mapped into a 1-dimensional bounding box, denoted by $\gamma(P^{\text{attr, opt, val, wt}}(x)) = \mathcal{B}^1$, as shown in Table I, where ϵ is the machine's epsilon. Therefore, a predicate $P(x)$ is covered by $\mathcal{B}^1 = [\min_1, \max_1]$ only if the permitted values defined by $P(x)$ lie in $[\min_1, \max_1]$; for brevity, we say the predicate is enclosed by \mathcal{B}^1 . Similarly, we say an expression Ω over \mathbb{R}^n is partially enclosed by \mathcal{B}^h w.r.t. the projection $\pi_{d_1 \dots d_h}$, denoted by $\Gamma^\pi(\Omega)$.

$$\Gamma^\pi(\Omega) = \{\mathcal{B}^h \mid \forall P_i^{\text{attr, opt, val, wt}}(x) \in \pi(\Omega),$$

$$\gamma(P_i(x)) \cap \xi_{P_i^{\text{attr}}}(\mathcal{B}^h) = \gamma(P_i(x))\}.$$

Furthermore, we say an expression Ω is fully enclosed by \mathcal{B}^k when the identity projection $\pi_{\mathcal{I}}$ is applied, $\Gamma^{\pi_{\mathcal{I}}}(\Omega) = \mathcal{B}^k$. Lastly, the smallest \mathcal{B}^k , or the minimum bounding box (MBB), that (partially) encloses an expression Ω is given by

$$\Gamma_{\min}^\pi(\Omega) = \{\arg \min_{\mathcal{B}_i^k} \Gamma^\pi(\Omega) = \mathcal{B}_i^k\}.$$

The Boolean expression Ω is said to have size k , denoted by $|\Omega| = k$, when having k predicates; hence, Ω is represented by $\Gamma_{\min}^\pi(\Omega)$ defined over a k -dimensional subspace.

¹The mapping strategy for predicates with operator \in, \notin, \neq is especially effective because the number of predicates per expression is on the order of tens while the number of space dimensions is on the order of thousands.

C. Top-k Matching Semantics

Our formulation of subscriptions and events as expressions enables us to support a wide range of matching semantics, including the classical pub/sub matching: *Given an event ω and a set of subscriptions Ω , find all subscriptions $\Omega_i \in \Omega$ that are satisfied by ω .* We refer to this problem as the *stabbing subscription*² $\text{SQ}(\omega)$ and formalize it as follows:

$$\text{SQ}(\omega) = \{\Omega_i \mid \forall P_q^{\text{attr, opt, val, wt}}(x) \in \Omega_i, \exists P_o^{\text{attr, opt, val, wt}}(x) \in \omega, P_q^{\text{attr}} = P_o^{\text{attr}}, \exists x \in \text{Dom}(P_q^{\text{attr}}), P_q(x) \wedge P_o(x)\}.$$

Alternatively, we can (approximately) express *stabbing subscription* as subspace matching in \mathbb{R}^n as follows, where the approximation is due to the mapping function γ :

$$\text{SQ}(\omega) = \{\Omega_i \mid \forall P_q^{\text{attr, opt, val, wt}}(x) \in \Omega_i, \exists P_o^{\text{attr, opt, val, wt}}(x) \in \omega, P_q^{\text{attr}} = P_o^{\text{attr}}, \gamma(P_q(x)) \cap \gamma(P_o(x)) \neq \emptyset\}.$$

Finally, we adopt the popular vector space scoring used in information retrieval (IR) systems³ for computing the score of a matched subscription Ω_i for a given event ω (to enable top-k computation), denoted by $\text{score}(\omega, \Omega_i)$, and defined by

$$\text{score}(\omega, \Omega_i) = \sum_{P_q(x) \in \Omega_i, P_o(x) \in \omega, P_q^{\text{attr}} = P_o^{\text{attr}}} P_q^{\text{wt}} \times P_o^{\text{wt}}.$$

Similarly, we compute the upper bound score for an event ω w.r.t. an upper bound weight-summary ($\text{sum}_{\text{wt}}^\Omega$) for a set of subscriptions Ω as follows

$$\text{uscore}(\omega, \text{sum}_{\text{wt}}^\Omega) = \sum_{P_o(x) \in \omega} P_o^{\text{wt}} \times \text{sum}_{\text{wt}}^\Omega(P_o^{\text{attr}}),$$

where $\text{sum}_{\text{wt}}^\Omega(\text{attr})$ returns the upper bound score of attr over the set of subscriptions Ω which is given by

$$\text{sum}_{\text{wt}}^\Omega(\text{attr}) = \max_{\Omega_i \in \Omega, P_q(x) \in \Omega_i, P_q^{\text{attr}} = \text{attr}} P_q^{\text{wt}}.$$

IV. BE*-TREE STRUCTURAL PROPERTIES

BE*-Tree is generic index structure for indexing a large collection of Boolean expressions (i.e., subscriptions) and for efficient retrieval of the most relevant matching expressions given a stream of incoming expressions (i.e., events). BE*-Tree supports Boolean expressions with an expressive set of operators defined over a high-dimensional continuous space. BE*-Tree copes with the curse of dimensionality challenge through a non-rigid two-phase space-cutting technique that significantly reduces the complexity and the level of uncertainty of choosing an effective criterion to recursively cut the space and that identifies highly dense subspaces. The two phases BE*-Tree employs are: (1) *partitioning* which is the global structuring to determine the next best attribute attr_i (i.e., the i^{th} dimension in \mathbb{R}^n) for splitting the space and (2) *non-rigid clustering* which is the local structuring for each partition to determine the best grouping of expressions w.r.t. the expressions' range of values for attr_i . BE*-Tree not only supports dynamic insertion and deletion of expressions, but it

²This is a relaxation of the stabbing query problem, in which interval cutting is generalized to subspace cutting.

³We are not limited to IR scoring, but support any monotonic scoring function.

also adapts to workload changes by incorporating top-down (data and space clustering) and bottom-up (space clustering) expansion within each clustering phase.

The data clustering aims to avoid indexing empty space and to adapt to a skewed workload while space clustering aims to avoid degeneration of the structure in the presence of frequent insertions and deletions of expressions. Conceptually, the space and the data clustering techniques are a hybrid scheme that takes the best of both worlds. On the one hand, the data clustering employs data dependent grouping of expressions to adapt to different workload distributions and to avoid indexing empty space, and on the other hand, the space clustering employs space dependent grouping of expressions to accommodate an insertion-independent mechanism.

In general, BE*-Tree is an n -ary tree structure in which a leaf node contains the actual data (expressions) and an internal node contains partial information about data (e.g., an attribute and a range of values) in its descendant leaf nodes. BE*-Tree consists of three classes of nodes: p -node (partition node) for storing the partitioning information, i.e., an attribute; c -node (cluster node) for storing the clustering information, i.e., a range of values; and l -node (leaf node), being at the lowest level of the tree, for storing the actual data. Moreover, p -nodes and c -nodes are logically organized in a special directory structure for fast tree traversal and search space pruning. Thus, a set of p -nodes are organized in a p -directory (partition directory), and a set of c -nodes are organized in a p -directory (cluster directory). The overall BE*-Tree structure together with its structural properties is depicted in Figure 2.

A. Non-rigid Space-cutting

BE*-Tree's non-rigid two-phase space-cutting, the partitioning followed by the non-rigid clustering, introduces new challenges such as how to determine the right balance between the partitioning and clustering, and how to develop a robust principle to alternate between both. In BE*-Tree, we propose a splitting policy to guide the non-rigid clustering phase for establishing not only a robust principle for alternating between partitioning and clustering but also for naturally adapting to the workload distributions. We begin discussing the overall structure and the two-phase cutting dynamics of BE*-Tree before presenting the key design principles behind BE*-Tree.

In BE*-Tree, the partitioning phase, conceptually a global adjusting mechanism, is the first phase of our space-cutting technique. The partitioning phase is invoked once a leaf node overflows (initially occurs at the root level). This phase involves ranking each candidate attr_i , attributes that appear in the expressions of the overflowed l -node, in order to determine the most effective attr_i for partitioning, as outlined in Algorithm 4. Essentially, this process identifies the highest ranking attr_i , given a scoring function (Section V-B), to spread expressions into smaller groups (leaf nodes). In short, partitioning space based on a high-ranking attr_i enables the pruning of search space more efficiently while coping with the curse of dimensionality by considering a single attr_i for each partitioning phase.

Upon executing the partitioning phase, the high-dimensional indexing problem is reduced to one-dimensional interval in-

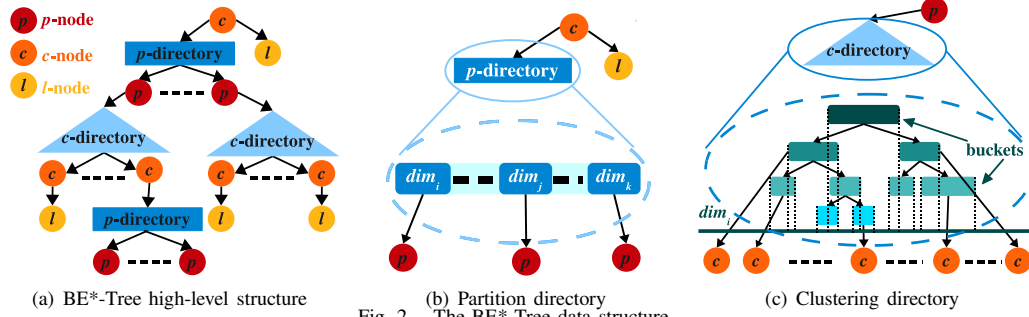


Fig. 2. The BE*-Tree data structure.

dexing, which paves the way to exploit underlying distribution of a single attr_i at a time through BE*-Tree's clustering phase, conceptually a local adjusting mechanism. At the core of our clustering phase, and ultimately of our BE*-Tree, are (1) a non-rigid clustering policy to group overlapping expressions (into buckets) that minimizes the overlap and the coverage among these buckets and (2) a robust and well-defined policy to alternate between the partitioning and the clustering (summarized in Algorithm 5).

The absence of a well-defined policy gives rise to the dilemma of whether to further pursue the space clustering or to switch back to the partitioning. Furthermore, a partitioned bucket can no longer be split without suffering from the cascading split problem [14]. Thus, a clustering policy that cannot react and adapt to the insertion sequence is either prone to ineffective buckets that do not take advantage of the domain selectivity for effectively pruning the search space or is prone to suffering from substantial performance overhead due to the cascading split problem. Therefore, a practical space clustering must support dynamic insertion and deletion and must adapt to any workload distributions, yet satisfying the cascading-split-free property.

In the clustering phase, each group of expressions is referred to as a bucket. Formally, a *bucket* is a 1-dimensional bounding box over attr_i , denoted by \mathcal{B}^1 , and an expression Ω with predicate $P_j^{\text{attr}, \text{opt}, \text{val}, \text{wt}}(x) \in \Omega$ is assigned to a bucket only if

$$\text{attr}_i = P_j^{\text{attr}} \quad \text{and} \quad \gamma(P_j(x)) \cap \xi_{P_j^{\text{attr}}}(\mathcal{B}^1) = \gamma(P_j(x)).$$

Furthermore, a bucket has a minimum bounding box (MBB) that partially encloses all of its expressions Ω if and only if

$$\mathcal{B}^1 = \bigcup_{\Omega_j \in \Omega} \Gamma_{\min}^{\pi_{\text{attr}_i}}(\Omega_j).$$

Moreover, each bucket is associated with exactly one c -node in the BE*-Tree, which is responsible for storing and maintaining information about the bucket's assigned expressions.

B. Structural Adaptation

The problem of supporting dynamic insertion and deletion is further complicated due to the non-rigid and dynamic clustering introduced by BE*-Tree. In a rigid clustering that relies on deterministically dividing the space in half, a best-effort deletion strategy is adequate because the size of the clustering boundary is deterministic and fixed; thereby, removing an expression does not alter the boundary. However, the non-rigid clustering directory of BE*-Tree defines a relative

(non-deterministic) bucket boundary as minimum boundary to enclose all of the bucket's expressions. Consequently, a simple best-effort approach for maintaining the boundary is inadequate because for a pathological case the boundary in BE*-Tree can unnecessarily span a large space, if an expression, defined over a large subspace, is inserted and later removed. To overcome this bucket boundary degeneration, BE*-Tree introduces a set of inactive dimensions which are used exclusively for efficiently keeping the boundary size up-to-date in presence of insertions and deletions. To incorporate the concept of inactive dimensions, we revise the bucket definition as follows:

A *bucket* at the k^{th} level of BE*-Tree is a k -dimensional bounding box, \mathcal{B}^k over \mathbb{R}^k , where $\pi_{d_1 \dots d_k}(\mathbb{R}^n) = \mathbb{R}^k$, and each $d_i \in \{d_1 \dots d_k\}$ is chosen by the partitioning phase at the i^{th} level of BE*-Tree. Also, we refer to the first $k-1$ dimensions of \mathcal{B}^k as inactive while we refer to the k^{th} dimension (where $d_k \in \mathbb{R}^n$) as active dimension. Furthermore, the bucket's partial enclosure property w.r.t. its expressions Ω holds if and only if

$$\mathcal{B}^k = \bigcup_{\Omega_j \in \Omega} \Gamma_{\min}^{\pi_{d_1 \dots d_k}}(\Omega_j).$$

The $k-1$ inactive dimensions correspond to $k-1$ attributes selected by the partitioning phase on the path from the root to the current clustering directory at the k^{th} -level of BE*-Tree. These inactive dimensions are utilized only during deletion while the active dimension is used for matching, insertion, and deletion. Without the notion of inactive dimension, deleting an expression results in the *cascading update problem*. The cascading update problem is an unpredictable chain reaction that propagates downwards in order to compute a new bucket boundary after an expression in one of the subtree's buckets has been removed. Thus, at every level of the tree, a brute force boundary re-computation of every bucket is required. However, with inactive dimensions in place, the boundary re-computation is required only for the bucket that is associated with the leaf node from which the expression was removed (or updated). Moreover, this information needs to be only sent upward along the path to the root. Thus, at the higher levels of BE*-Tree, the new boundary can simply be reconstructed by considering the updated boundary information together with the boundary of unmodified buckets. The maintenance of active and inactive dimensions is handled by Algorithm 3.

C. BE*-Tree Invariance

Before stating BE*-Tree's invariance and operational semantics, we must distinguish among four bucket types: *open*

bucket: a bucket that is not yet partitioned; *leaf bucket*: a bucket that has no children (or has not been split); *atomic bucket*: a bucket that is a single-valued bucket which cannot further be split; and *home bucket*: a bucket that is the smallest existing bucket that encloses the inserting expression.

The correctness of the BE*-Tree operational semantics is achieved based on the following three rules:

- 1) *insertion rule*: an expression is always inserted into the smallest bucket that encloses it, while respecting the descendant-repelling property (Section V-A3).
- 2) *forced split rule*: an overflowing non-leaf bucket is always split before switching back to the partitioning.
- 3) *merge rule*: an underflowing leaf bucket is merged with its parent only if the parent is an open bucket

Finally, the BE*-Tree correctness can be summarized as follows: INVARIANCE: *Every expression Ω is always inserted into the smallest bucket that encloses it and a non-atomic bucket is always split first before it is partitioned.*

V. BE*-TREE ADAPTIVENESS PROPERTIES

An essential property of any index structure is to dynamically adapt as the workload changes in order to prevent index deterioration. Supporting adaptiveness is central to the design of BE*-Tree and has shaped every aspect of BE*-Tree. The adaptive nature of BE*-Tree is reflected in its bi-direction expansion that foresees index evolution, in its splitting strategy that improves search-space pruning w.r.t. the subscription workload, and in its proposed simplified ranking function that actively refines the space partitioning based on both subscription and event workloads. The core of BE*-Tree adaptiveness is captured by its insertion operation. Also, the adaptiveness in deletion is accomplished through maintenance of inactive dimension, which was discussed in Section IV-B, and through recycling nodes [7] (i.e., reinsertion policy) whose rank have dropped below certain threshold, our ranking function is defined in Section V-B.

A. Bi-directional Expansions

The main aim of the bi-directional expansion is to determine the right balance between a rigid clustering of the grid-based approach (space dependent) and a flexible and a dynamic clustering similar to that of an *R*-Tree-based approach (data dependent), while avoiding the cascading split problem. Therefore, the final objective is to enable BE*-Tree to model any data distribution (e.g., skewed), and only in the worst case degenerate to a grid-based approach which is best suited for a uniform distribution.

1) *Top-down Expansion*: The top-down expansion is part of the clustering phase in which an overflowing bucket is recursively split while keeping the bucket capacity within a threshold, denoted by A . There are two main types of bucket splitting: a data dependent splitting strategy (data clustering), discussed in Section V-A3, and a space dependent splitting strategy (space clustering) which follows a grid-based splitting, also in Section V-A3. As a result, whenever the data clustering approach fails to split the data into two or more buckets, while satisfying a minimum capacity requirement, denoted by α , for each bucket, it relies on space clustering

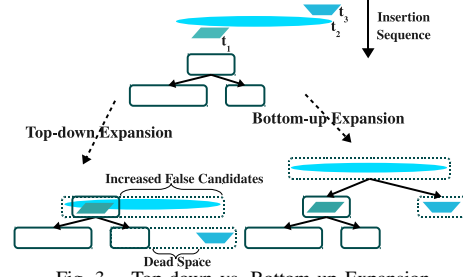


Fig. 3. Top-down vs. Bottom-up Expansion

such that the bucket is split in half, as shown in Algorithm 2, Lines 11-25.

2) *Bottom-up Expansion*: The bottom-up expansion reverses the conventional tree top-down expansion. It focuses on dynamically adjusting BE*-Tree's clustering directory to accommodate workload changes. Most importantly, the bottom-up expansion avoids BE*-Tree degeneration by preventing the partitioned bucket's size to expand indefinitely. Therefore, to achieve this goal, the clustering directory is re-adjusted by injecting a new bucket in the upward direction, whenever, a partitioned bucket is required to host an expression which is β times larger than the size of all expressions in its subtree. This procedure is illustrated in Figure 3 and its pseudocode is provided in Algorithm 2, Lines 6-10.

3) *Descendant-repelling & Overlap-free Split*: We present a novel splitting strategy for overflowing buckets, a strategy that falls in between the traditional overlapping-style splitting and grid-style rigid splitting. Our new splitting algorithm aims at systematically reducing the overlap, similar to grid-based structures, while adapting to the data distribution, similar to *R*-Tree-based structures. To achieve these properties, the overflowing bucket is split into three buckets: left, overlapping, and right buckets, such that the left and right buckets (descendant buckets) repel each other. Although their boundary can be dynamically changed to adapt to the data distribution, they can never overlap. Consequently, data that spans over the boundary of the descendant buckets is placed in the overlapping node. Thus, by definition the overlapping bucket is the home bucket to all of its data, and upon overflowing, the overlapping bucket is handled by switching back to the partitioning phase.

An important property of the overlap-free splitting, guaranteed under the uniform distribution, is realized when a *point event* (i.e., consisting of only equality predicates) falls within the overlapping region. This results in a substantial reduction of the search space in which only $\frac{2}{3}$ of the data must be scanned as opposed to checking all of the data when only two buckets are used and bucket-overlap is allowed. The $\frac{2}{3}$ -reduction in search space is due to the proposed splitting of overflowing buckets into three buckets, a left child, a right child, and an overlapping bucket such that left and right children never intersect, which is enforced by our descendant-repelling property. This overlap-free splitting guarantees that a point-event can cut through at most one child. Thus, in the worst case only $\frac{2}{3}$ of the data is visited⁴. The resulting reduc-

⁴The $\frac{2}{3}$ -reduction can simply be applied to extended events (events with non-equality predicates), if the descendant-repelling property is generalized to enforce that sibling buckets must be separated by the average extension size of events, where the extension is defined as $\gamma(\mathcal{P}^{\text{attr}, \text{opt}, \text{val}, \text{wt}}(x))$.

tion is substantial because an entire BE*-Tree subtree could be associate with each bucket. Finally, under the descendant-repelling property, all subscriptions under any sibling buckets are guaranteed to be disjoint on at least one dimension. Next, we present two algorithms to construct descendant-repelling and overlap-free splitting.

Exhaustive Splitting Algorithm Our first algorithm finds an effective splitting hyperplane to divide overflowing buckets into three buckets. In the splitting of the overflowing bucket within the clustering phase at the k^{th} level of BE*-Tree, first all the expressions $\Omega_i \in \Omega$ in the overflowing bucket are sorted w.r.t. the active dimension (k^{th} dimension of the overflowing bucket, denoted by attr^k) of BE*-Tree based on both upper and lower boundary of $P^{\text{attr}}(x) \in \Omega_i$, where $P^{\text{attr}} = \text{attr}^k$. The cost of the exhaustive splitting algorithm is dominated by the sorting phase together with enumeration over all possible splitting choices, namely, $O((n-2\alpha)n \log(n))$, where $n = |\Omega|$ and α is the minimum capacity of the resulting right and left buckets after splitting. In brief, the algorithm first assigns the I leftmost expressions to the left bucket (or I rightmost expressions to the right bucket); second, it distributes the remaining expressions to the right (or the left) and the overlapping buckets accordingly. The main objective of the algorithm is to find the optimal I that minimizes overlap and coverage while satisfying the descendant-repelling property and minimum bucket capacity α . The splitting is formally defined as follow:

$$\text{buckets}_{\text{opt}} = \begin{cases} \text{split}_{\text{dataClustering}}^{\text{opt}}(\text{find}_{\text{opt}}(\Omega_{\text{low}}, \Omega_{\text{up}})) & \text{if } \frac{|\mathcal{B}_R^k|}{|\mathcal{B}_L^k|} \geq \alpha \wedge \frac{|\mathcal{B}_O^k|}{|\mathcal{B}_L^k|} \geq \alpha \\ \text{split}_{\text{spaceClustering}}(\Omega, s) & \text{otherwise,} \end{cases}$$

where $s = \chi_k(\mathcal{B}^k)$, \mathcal{B}^k is the overflowing bucket at the k^{th} level of BE*-Tree, and $\mathcal{B}_R^k, \mathcal{B}_O^k, \mathcal{B}_L^k$ are the new right, overlapping, and left buckets, respectively, returned by either split functions. Furthermore, the most cost effective splitting, denoted by find_{opt} , is determined by a set of buckets (right, overlapping, left) which minimizes the total overlap, to break ties among candidate bucket sets, the set with the smallest coverage is selected. The function find_{opt} is formally defined as follows:

$$(\Omega_j, I) \leftarrow \text{find}_{\text{opt}}(\Omega_{\text{low}}, \Omega_{\text{up}}) \leftarrow \arg \min_{(\Omega_j, I) \in \text{min}_{\text{overlap}}(\Omega_{\text{low}}, \Omega_{\text{up}})} \text{coverage}(\Omega_j, I)$$

$$\{(\Omega_j, I)\} \leftarrow \text{min}_{\text{overlap}}(\Omega_{\text{low}}, \Omega_{\text{up}}) \leftarrow \left\{ \arg \min_{\substack{\Omega_j \in \{\Omega_{\text{low}}, \Omega_{\text{up}}\}, \\ I \in (\alpha \cdots |\Omega| - \alpha)}} \text{overlap}(\Omega_j, I) \right\},$$

where Ω_{low} represents the expressions in the overflowing bucket \mathcal{B}^k sorted in descending order w.r.t. the lower bound of the expressions' k^{th} dimension; similarly, Ω_{up} represents the expressions sorted in ascending order w.r.t. the upper bound on the k^{th} dimension.

The overlap and coverage are computed over k^{th} dimension of the bucket as follows:

$$\text{overlap}(\Omega, I) = \left(\lambda_k(\pi_k(\mathcal{B}_O^k)) \mid \mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k \leftarrow \text{split}_{\text{dataClustering}}^{\text{opt}}(\Omega, I) \right)$$

$$\text{coverage}(\Omega, I) = \left(\lambda_k(\pi_k(\mathcal{B}_L^k) \cup \pi_k(\mathcal{B}_O^k) \cup \pi_k(\mathcal{B}_R^k)) \mid \mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k \leftarrow \text{split}_{\text{dataClustering}}^{\text{opt}}(\Omega, I) \right).$$

The descendant-repelling and overlap-free buckets are defined as follows⁵, which is the basis of our data clustering procedure:

$$(\mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k) \leftarrow \text{split}_{\text{dataClustering}}^{\text{opt}}(\Omega, I) = \bigcup_{i=1}^I \left(\mathcal{B}_L^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_i) \right), \\ \bigcup_{i=I+1}^{|\Omega|} \left((\mathcal{B}_R^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_i) \mid \pi_k(\mathcal{B}_L^k) \cap \gamma(\pi_{d_k}(\Omega_i)) = \emptyset), \right. \\ \left. (\mathcal{B}_O^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_i) \mid \pi_k(\mathcal{B}_L^k) \cap \gamma(\pi_{d_k}(\Omega_i)) \neq \emptyset) \right).$$

In contrast, if no data clustering exists that produces left and right buckets with at least α expressions, then we rely on space clustering, a default clustering which divides the space in half in order to avoid degeneration of the clustering directory.

$$(\mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k) \leftarrow \text{split}_{\text{spaceClustering}}(\Omega, s) = \bigcup_{i=1}^{|\Omega|} \left((\mathcal{B}_L^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_i) \mid \mu(\gamma(\pi_{d_k}(\Omega_i))) > s), \right. \\ \left. (\mathcal{B}_R^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_i) \mid M(\gamma(\pi_{d_k}(\Omega_i))) < s), \right. \\ \left. (\mathcal{B}_O^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_i) \mid s \cap \gamma(\pi_{d_k}(\Omega_i)) \neq \emptyset) \right).$$

Approximate Splitting Algorithm Unlike the exhaustive algorithm, the cost of our fast splitting algorithm is dominated by sorting expressions, namely, $O(n \log(n))$, where $n = |\Omega|$. In a nutshell, the approximate algorithm assigns expressions one-by-one from the leftmost expression to the left bucket and from the rightmost expression to the right bucket while pushing expressions that span both left and right bucket into the overlap bucket. Formally the algorithm is defined as follows:

$$\text{buckets}_{\text{apx}} = \begin{cases} \text{split}_{\text{dataClustering}}^{\text{apx}}(\Omega_{\text{low}}, \Omega_{\text{up}}) & \text{if } \frac{|\mathcal{B}_R^k|}{|\mathcal{B}_L^k|} \geq \alpha \wedge \frac{|\mathcal{B}_O^k|}{|\mathcal{B}_L^k|} \geq \alpha \\ \text{split}_{\text{spaceClustering}}(\Omega, s) & \text{otherwise,} \end{cases}$$

where $s = \chi_k(\mathcal{B}^k)$.

$$(\mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k) \leftarrow \text{split}_{\text{dataClustering}}^{\text{apx}}(\Omega_{\text{up}}, \Omega_{\text{low}}) = \bigcup_{i=1}^{|\Omega_{\text{up}}|} \left((\mathcal{B}_L^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_{\text{up}_i}) \mid \gamma(\pi_{d_k}(\Omega_{\text{up}_i})) \cap \pi_k(\mathcal{B}_R^k) = \emptyset), \right. \\ \left. (\mathcal{B}_R^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_{\text{low}_i}) \mid \gamma(\pi_{d_k}(\Omega_{\text{low}_i})) \cap \pi_k(\mathcal{B}_L^k) = \emptyset), \right. \\ \left. \bigcup_{i=1}^{|\Omega_{\text{up}}|} (\mathcal{B}_O^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \cdots \pi_{d_k}}(\Omega_{\text{up}_i}) \mid \Omega_{\text{up}_i} \notin \mathcal{B}_L^k \wedge \Omega_{\text{up}_i} \notin \mathcal{B}_R^k) \right).$$

⁵The provided definition of $\text{split}_{\text{dataClustering}}$ is used when computing buckets over Ω_{up} ; for applying it to Ω_{low} , the computation of \mathcal{B}_L^k and \mathcal{B}_R^k must simply be reversed.

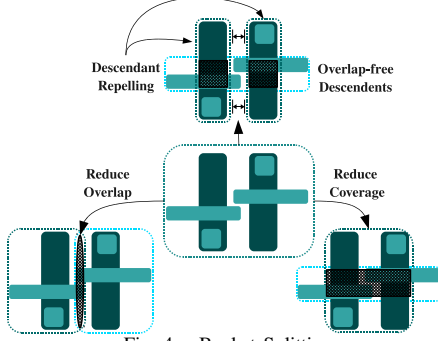


Fig. 4. Bucket Splitting

B. Adaptive Ranking Function

Our ranking function simplifies BE-Tree's ranking function while maintaining the actual index objective, namely, reducing the matching cost. The matching cost is defined as the total number of predicate evaluations: *minimizing false candidate computations* and *minimizing true candidate computations*. The former cost is the total number of predicate evaluations for each discarded search path or for each unsatisfied expression, i.e., penalizing multiple search paths or paths resulting in many false positives. The latter cost is the number of predicate evaluations before concluding a set of expressions as matched, i.e., promoting evaluation of common predicates among expressions exactly once.

Furthermore, our ranking model relies on the notion of *false candidates*: the scanned expressions (subscriptions) that are not matched by an input expression (events). Our model assigns a rank to each node n_i using the function $\text{rank}(n_i)$ which is a combination of the $\text{loss}(n_i)$ and $\text{gain}(l_j)$ functions. $\text{loss}(n_i)$ computes for each node the false candidates generated over a window of m events. $\text{gain}(l_j)$ computes the number of common predicates for each of its expressions. Formally, $\text{rank}(n_i)$, $\text{loss}(n_i)$, and $\text{gain}(l_j)$ are defined as follows:

$$\text{rank}(n_i) = \begin{cases} \text{gain}(n_i) - \text{loss}(n_i) & \text{if } n_i \text{ is a } l\text{-node} \\ \left(\sum_{n_j \in \text{des}(n_i)} \text{rank}(n_j) \right) - \text{loss}(n_i) & \text{otherwise,} \end{cases}$$

where $\text{des}(n_i)$ returns n_i 's immediate descendants

$$\text{loss}(n_i) = \sum_{e' \in \text{window}_m(n_i)} \frac{\# \text{ discarded pred eval for } e'}{|\text{window}_m(n_i)|}$$

$$\text{gain}(l_j) = \# \text{ common pred.}$$

The proposed ranking model is simply generalized for splitting an overflowing node l_j -node using a new attr_i . It is given by

$$\text{rank}(l_i) = \text{gain}(l_i) - \text{loss}(l_i),$$

where $\text{gain}(l_i)$ is approximated by the number of expressions that have a predicate on attr_i , and $\text{loss}(l_i)$ is estimated either using the selectivity by constructing a histogram, in which $\text{loss}(l_i)$ is the average bucket size in the histogram, or using popularity, an optimistic approach, by initially setting $\text{loss}(l_i)$ to zero to eliminate any histogram construction and to rely on

Algorithm 1 Insert(Ω , $cnode$, $cdir$)

```

1: /* Find attr with max rank not yet used for partitioning */
2: if  $cnode.pdir \neq \text{NULL}$  then
3:   for  $P_i^{\text{attr}} \in \omega$  do
4:     if !IsUsed( $P_i^{\text{attr}}$ ) then
5:        $pnode \leftarrow \text{SearchPDir}(P_i^{\text{attr}}, cnode.pdir)$ 
6:       if  $pnode \neq \text{NULL}$  then
7:          $\text{foundPartition} = \text{true}$ ;
8:         if  $\text{maxRank} < pnode.rank$  then
9:            $\text{maxPnode} \leftarrow pnode$ 
10:           $\text{maxRank} \leftarrow pnode.rank$ 
11: /* if no partitioning found, then insert into the l-node */
12: if ! $\text{foundPartition}$  then
13:   Insert( $\Omega$ ,  $cnode.lnode$ )
14:   UpdateLNodeTopNScore( $\Omega$ ,  $cnode.lnode.sum_{\text{nt}}^{\Omega}$ )
15: /* if c-node is the root, then do partitioning; otherwise clustering */
16: if isRoot( $cnode$ ) then
17:   Partitioning( $cnode$ )
18: else
19:   Clustering( $cdir$ )
19:  $\text{maxCdir} \leftarrow \text{InsertCDir}(\Omega, \text{maxPnode.cdir})$ 
20: insert( $\Omega$ ,  $\text{maxCdir.cnode}$ ,  $\text{maxCdir}$ )
21: UpdatePNodeScore( $\text{maxPnode}$ )
22: UpdatePNodeTopNScore( $\Omega$ ,  $\text{maxPnode.sum}_{\text{nt}}^{\Omega}$ )

```

Algorithm 2 InsertCDir(Ω , $cdir$)

```

1:  $\mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k \leftarrow cdir.\mathcal{B}_L^k, cdir.\mathcal{B}_O^k, cdir.\mathcal{B}_R^k$ 
2:  $\mathcal{B}^k \leftarrow \mathcal{B}_L^k \cup \mathcal{B}_O^k \cup \mathcal{B}_R^k$ 
3: if IsLeaf( $\mathcal{B}_O^k$ ) then
4:   UpdateBucketInfo( $\Omega, \mathcal{B}_O^k$ )
5:   return  $cdir$ 
6: /* Bottom-up: if total  $k^{\text{th}}$  boundary is expanded by more than  $\beta$  */
7: if  $\lambda_k(\gamma(\pi_{d_k}(\Omega)) \cup \mathcal{B}^k) > \beta \times \lambda_k(\pi_k(\mathcal{B}^k))$  then
8:    $\mathcal{B}_{\text{new}}^k \leftarrow \Omega$ 
9:    $cdir \xleftarrow{\text{insert}} \mathcal{B}_{\text{new}}^k \xleftarrow{\text{new parent}} \mathcal{B}_L^k, \mathcal{B}_O^k, \mathcal{B}_R^k$ 
10:  UpdateBucketInfo( $\Omega, \mathcal{B}_{\text{new}}^k$ )
11:  return  $cdir$ 
12: /* Top-down: otherwise */
13: else
14:   /* If descendant-repelling property is violated by inserting  $\Omega$  */
15:   if  $\gamma(\pi_{d_k}(\Omega)) \cap \pi_k(\mathcal{B}_L^k) \neq \emptyset \wedge \gamma(\pi_{d_k}(\Omega)) \cap \pi_k(\mathcal{B}_R^k) \neq \emptyset$  then
16:     UpdateBucketInfo( $\Omega, \mathcal{B}_O^k$ )
17:     return  $cdir$ 
18:   else if  $\gamma(\pi_{d_k}(\Omega)) \cap \pi_k(\mathcal{B}_L^k) \neq \emptyset$  then
19:     return InsertCDir( $\Omega, \mathcal{B}_L^k$ )
20:   else if  $\gamma(\pi_{d_k}(\Omega)) \cap \pi_k(\mathcal{B}_R^k) \neq \emptyset$  then
21:     return InsertCDir( $\Omega, \mathcal{B}_R^k$ )
22:   else
23:     /* insert into a bucket which results in minimum expansion */
24:     if  $\lambda_k(\gamma(\pi_{d_k}(\Omega)) \cup \mathcal{B}_L^k) < \lambda_k(\gamma(\pi_{d_k}(\Omega)) \cup \mathcal{B}_R^k)$  then
25:       return InsertCDir( $\Omega, \mathcal{B}_L^k$ )
26:     else
27:       return InsertCDir( $\Omega, \mathcal{B}_R^k$ )

```

the matching feedback mechanism for adjusting the ranking if necessary. In our experimental evaluation, the optimistic approach results in an improved matching and insertion time.

Algorithm 3 UpdateBucketInfo(Ω , $cdir.B^k$)

```

1: /* Update the active dimension (the bucket  $k^{th}$  dimension) */
2:  $B^k \leftarrow \gamma(\pi_{d_k}(\Omega)) \cup \pi_k(B^k)$ 
   /* Update the inactive dimensions  $(1 \dots (k-1)^{th}$  dimensions) */
3:  $B^k \leftarrow \Gamma_{\min}^{\pi_{d_1} \dots \pi_{d_{k-1}}}(\Omega) \cup \pi_{1 \dots (k-1)}(B^k)$ 
4: UpdateCNodeInactiveTopNScore( $\Omega_i, B^k.cnode.sum_{wt}^{\Omega}$ )

```

Algorithm 4 Partitioning($cnode$)

```

1:  $lnode \leftarrow cnode.lnode$ 
2: while IsOverflowed( $lnode$ ) do
3:    $attr \leftarrow \text{GetNextHighestScoreUnusedAttr}(lnode)$ 
   /* Create new partition for the next highest ranked attr */
4:    $pnode \leftarrow \text{CreatePDir}(attr, cnode.pdir)$ 
   /* Move all the subscriptions with predicate on attr */
5:   for  $\Omega_i \in lnode.\Omega$  do
6:     if  $\exists P_j^{attr} \in \Omega_i, P_j^{attr} = attr$  then
7:        $cdir \leftarrow \text{InsertCDir}(\Omega_i, pnode.cdir)$ 
8:        $\text{Move}(\Omega_i, lnode, cdir.cnode.lnode)$ 
9:        $\text{UpdatePNodeTopNScore}(\Omega_i, pnode.sum_{wt}^{\Omega})$ 
10:   $\text{Clustering}(pnode.cdir)$ 
11:  $\text{UpdateClusterCapacity}(lnode)$ 

```

Algorithm 5 Clustering($cdir$)

```

1:  $lnode \leftarrow cdir.cnode.lnode$ 
2: if !isOverflowed( $lnode$ ) then
3:   return
4: if !IsLeaf( $cdir$ ) or IsPartitioned( $cdir$ ) then
5:   Partitioning( $cdir.cnode$ )
6: else
7:   if Exhaustive Splitting Algorithm then
8:      $cdir.B_L^k, cdir.B_O^k, cdir.B_R^k \leftarrow \text{buckets}_{opt}$ 
9:   else
10:     $cdir.B_L^k, cdir.B_O^k, cdir.B_R^k \leftarrow \text{buckets}_{spx}$ 
11:   for  $\Omega_i \in lnode.\Omega$  do
12:     if  $\gamma(\pi_{d_k}(\Omega_i)) \cap \pi_k(cdir.B_L^k) = \gamma(\pi_{d_k}(\Omega))$  then
13:        $\text{UpdateBucketInfo}(\Omega_i, cdir.B_L^k)$ 
14:        $\text{Move}(\Omega_i, lnode, cdir.B_L^k.cnode.lnode)$ 
15:     else if  $\gamma(\pi_{d_k}(\Omega_i)) \cap \pi_k(cdir.B_R^k) = \gamma(\pi_{d_k}(\Omega))$  then
16:        $\text{UpdateBucketInfo}(\Omega_i, cdir.B_R^k)$ 
17:        $\text{Move}(\Omega_i, lnode, cdir.B_R^k.cnode.lnode)$ 
18:   Partitioning( $cdir.B_O^k.cnode$ )
19:   Clustering( $cdir.B_L^k$ )
20:   Clustering( $cdir.B_R^k$ )
21:  $\text{UpdateClusterCapacity}(lnode)$ 

```

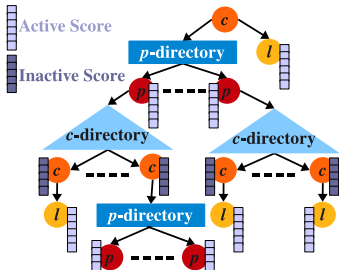


Fig. 5. BE*-Tree Top-k Structural Overview

VI. TOP-K PROCESSING

Our top-k algorithm is a two-stage process in which, on the one hand, a cost-based tree traversal is utilized to determine which p -node in BE*-Tree to visit next, and, on the other hand,

along the traversal path, a top-k upper bound computation determines whether an l -node content is scanned or not. For an efficient top-k implementation of the former stage, a p -node-max-heap structured is maintained, whereas for the latter, a k -min-heap structure is maintained.

The BE*-Tree's partitioning and clustering abstraction lends itself well for top-k processing. This abstraction also provides a different level of granularity for scoring. For instance, at a p -node level, in which only the information about the attribute space (and the upper bound score for each attribute) is maintained, no maintenance of scoring w.r.t. the distribution of values within each dimension is required. As a result, the upper bound scoring information that must be stored at each p -node can be compressed. This compression results in huge space reduction. Moreover, the same compression technique is applied for storing upper bound scores within each l -node. The overall BE*-Tree top-k structure is illustrated in Figure 5.

The top-k algorithm is as follows: Upon arrival of an event expression, the score between the event and the root's relevant p_i -nodes are computed to form a key ($\text{score}_{p_i}, p_i\text{-node}$) which is then inserted into the p -node-max-heap (Algorithm 6 Lines 6-12). Next, the p -node with highest score value is removed from the p -node-max-heap and its subtree is visited (Algorithm 6 Lines 13-19). For the chosen p -node, all of its relevant c_j -nodes (those c_j -nodes having a non-empty intersection with the event) are visited while going through the c -directory (Algorithms 7). For each c_j -node, its l -node is scanned. This process is applied recursively until either the top-k expressions are discovered or all the relevant portions of BE*-Tree are visited. We now discuss the pruning procedures (top-k early terminations):

- 1) **l -node pruning** Before scanning expressions in an l -node, given an event expression, the upper bound score of an l -node is computed, if the score is larger than the minimum score of the matched expressions seen so far, stored in k -min-heap, or there are fewer than k matched expressions seen so far, then the l -node is visited; it is pruned otherwise, as outlined in Algorithm 6, Line 2.
- 2) **p -node pruning** Before inserting the p -node into the p -node-max-heap, the score of p -node w.r.t. the event expression is computed; only if the score is larger than the minimum score of matched expressions (or fewer than k matched expressions are gathered so far), then the p -node is inserted. Similarly, before removing the p -node with the highest score from the p -node-max-heap, the score of the highest p -node is compared with the lowest score of matched expressions, if the p -node's score is smaller, then the entire matching process is terminated, resulting in a substantial saving in the matching time, as given in Algorithm 6, Lines 15-17 results.

VII. EVALUATIONS

We present an extensive evaluation of BE*-Tree based on synthetic and real datasets. The experiments were ran on a machine with two Intel Xeon 3.00GHz Quad-core processors having 16GB of memory and running CentOS 5.5 64bit. All algorithms are implemented in C and compiled using gcc 4.1.2.

Algorithm 6 MatchTopk($\omega, \Omega_{matched}, cnode$)

```

1:  $lnodeScore = uscore(\omega, lnode.sum_{wt}^{\Omega});$ 
2: if  $lnodeScore < MinScore(k-min-heap)$  then
3:   /* lnode pruning, skipping an l-node */
4: else
5:    $\Omega_{matched} \leftarrow FindMatchingExpr(cnode.lnode)$ 
6:   /* Iterate over event's predicates */
7:   for  $P_i^{attr} \in \omega$  do
8:     /* Retrieve c-node's p-directory for  $P_i^{attr}$  */
9:      $pnode \leftarrow SearchPDir(P_i^{attr}, cnode.pdir)$ 
10:    /* If  $P_i^{attr}$  exists in the p-directory */
11:    if  $pnode \neq NULL$  then
12:       $pnodeScore = uscore(\omega, pnode.sum_{wt}^{\Omega});$ 
13:      if  $pnodeScore > MinScore(k-min-heap)$  then
14:         $p-node-max-heap \leftarrow (pnode, pnodeScore)$ 
15:      while  $!IsEmpty(k-min-heap)$  do
16:         $maxScorePnode \leftarrow GetMax(p-node-max-heap)$ 
17:        if  $!IsFull(k-min-heap)$  then
18:          if  $maxScorePnode < MinScore(k-min-heap)$  then
19:             $break$  /* pnode pruning, matching termination */
20:         $SearchCDir(\omega, \Omega_{matched}, maxScorePnode.cdir)$ 
21:         $DeleteMax(p-node-max-heap)$ 

```

Algorithm 7 SearchCDir($\omega, \Omega_{matched}, cdir$)

```

1: if  $\gamma(\pi_{d_k}(\omega)) \cap \pi_k(cdir.B_O^k) \neq \emptyset$  then
2:    $MatchTopk(\omega, \Omega_{matched}, cdir.B_O^k, cnode)$ 
3: if  $\gamma(\pi_{d_k}(\omega)) \cap \pi_k(cdir.B_L^k) \neq \emptyset$  then
4:    $SearchCDir(\omega, \Omega_{matched}, cdir.B_L^k)$ 
5: if  $\gamma(\pi_{d_k}(\omega)) \cap \pi_k(cdir.B_R^k) \neq \emptyset$  then
6:    $SearchCDir(\omega, \Omega_{matched}, cdir.B_R^k)$ 

```

A. Experiment Overview

Our evaluation considers the following algorithms: (1) SCAN (a sequential scan), (2) k -ind (the conjunction algorithm implemented over k -index [2] in which all subscriptions are known in advance), (3) BE (BE-Tree [7]), (4) BE*-B (a batch-oriented construction of BE*-Tree in which all subscriptions are known in advance), (5) BE* (a fully dynamic version of BE*-Tree in which subscriptions are inserted individually). Most notably, BE*-Tree is the only structure that is designed to support continuous domains (with the exception of the naive SCAN). Therefore, to compare BE*-Tree with state-of-the-art algorithms, we emulate the effect of continuous domains using a large discrete domain. Furthermore, we primarily focus on the range operator (BETWEEN) over discrete domains. Moreover, we further restrict the size of ranges to only a handful of values, otherwise the enumeration technique of k -index [2] to support range operators would become infeasible due to exponential space explosion.

Our evaluation explores various crucial workload characteristics including *workload distribution*, *workload size*, *space dimensionality*, *average subscription and event size*, *dimension cardinality*, *number of clusters in each dimension*, *dimension cluster size*, and *event matching probability*. To generate various experimental configurations, we leverage the BEGen⁶ framework [7] including public domain data (i.e., DBLP repository data) to generate real workloads. Various workload

⁶<http://msrg.org/datasets/BEGen>

TABLE II
BE*-TREE MICRO EXPERIMENTS

(a) Max Leaf Capacity Trends		(b) Construction Time	
	A	Data Sets	Time (second)
Match Prob $< 1\%$	5	Unif (1M)	2.30
$1\% \leq$ Match Prob $< 10\%$	20	Zipf (1M)	2.45
Match Prob $\geq 10\%$	160	Author (760K)	4.98
		Title (250K)	2.83

configurations are captured in Table III. Due to limited space, in this paper, we primarily focus on a distinguishing subset of experimental results extracted from Table III.

B. Micro Experiments

BE*-Tree Parameters BE*-Tree internal parameters are extensively studied and the most robust configuration that resulted in the best overall performance in all experiments is presented here (the in-depth results are omitted). The values used throughout our experiments are A (maximum leaf capacity) shown in Table II(a); $\alpha = 10\%$ (minimum capacity of descendant buckets); $\beta = 4$ (maximum allowable bucket boundary increase before invoking bottom-up expansion); and the exhaustive splitting algorithm is used to handle overflowing buckets.

BE*-Tree Construction The dynamic construction time of BE*-Tree for our largest workload having up to a million subscriptions was below 5 seconds; BE*-Tree's average construction time over our main datasets are reported in Table II(b).

C. Macro Experiments

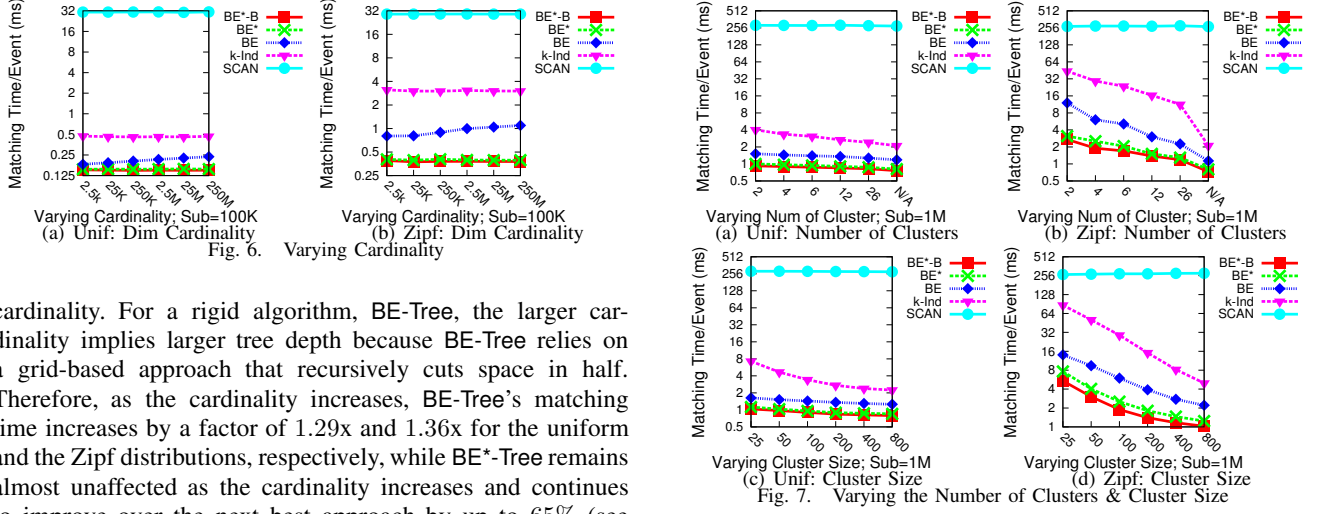
In our macro experiments, we first establish the effectiveness of BE*-Tree's internal techniques, namely, the non-rigid space-cutting, bi-directional tree expansion, and descendant-repelling and overlap-free splitting, then, we shift our focus to BE*-Tree's hierarchical top-k processing that is realized over these unique inner structures of BE*-Tree for finding the most relevant matches.

Workload Distribution A key distinguishing property that sets apart key- vs. non-key-based and rigid-based vs. non-rigid-based algorithms is the workload distribution. The effects of the distributions are shown in Figure 6-9 in which the graphs in the left column correspond to the uniform distribution (for choosing predicates' attributes) while the graphs on the right column correspond to the Zipf distribution with the exception of DBLP-driven workloads which follow a real-data distribution. The observed trend is that k -index, a non-key-based method, is highly sensitive to the underlying data distribution because few popular attributes appear in the majority of subscriptions which result in increased false candidates. The presence of popular attributes forces k -index to sort and scan through posting lists with large number of subscriptions [2]. The increased false candidate problem is reduced by BE-Tree, being a key-based approach. However, BE-Tree also needlessly indexes dead spaces. As a result, BE*-Tree on average outperforms k -index by 60% and 90% and BE-Tree by 35% and 65% for the uniform and the Zipf distributions, respectively.

Dimension Cardinality An important factor that further differentiates rigid from non-rigid algorithms is dimension

TABLE III
SYNTHETIC AND REAL WORKLOAD CONFIGURATIONS

	Workload Size	Number of Dimensions	Dimension Cardinality	Number of Clusters	Cluster Size	Sub/Event Size	Match Prob		DBLP (Author)	DBLP (Title)	Match Prob (Author)	Match Prob (Title)
Size	100K-1M	1M	100K	100K	100K	100K	1M		100-760K	50-250K	400k	150
Number of Dim	400	50-1400	400	400	400	400	400		677	677	677	677
Cardinality	250M	250M	2.5K-250M	250M	250M	250M	250M		26 ³	26 ³	26 ³	26 ³
Number of Clusters	4	4	4	2-26, N/A	4	4	4		—	—	4	4
Cluster Size	100	100	100	100	25-800	100	100		—	—	100	100
Avg. Sub Size	8	8	8	8	8	5-66	8		8	35	8	30
Avg. Event Size	15	15	15	15	15	13-81	15		8	35	16	43
Match Prob %	0.1	0.1	0.1	0.1	0.1	0.1	0.01-9		—	—	0.01-9	0.01-9



cardinality. For a rigid algorithm, BE-Tree, the larger cardinality implies larger tree depth because BE-Tree relies on a grid-based approach that recursively cuts space in half. Therefore, as the cardinality increases, BE-Tree's matching time increases by a factor of 1.29x and 1.36x for the uniform and the Zipf distributions, respectively, while BE*-Tree remains almost unaffected as the cardinality increases and continues to improve over the next best approach by up to 65% (see Figure 6a,b).

Another consequence of increasing the dimension cardinality is the effect on memory footprint. Especially for k -index, increasing dimension cardinality results in an exponential memory growth when the ratio of predicate range size to dimension cardinality is kept constant. The space blow-up occurs because k -index relies on the enumeration technique to support range predicates. For instance, in order to cope with the operator BETWEEN $[v_1, v_2]$, the enumeration essentially transforms the value of $v_2 - v_1$ from a decimal to a unary representation: an exponential transformation.

Number of Clusters/Cluster Size Next, we study the effects of the number of clusters and the cluster size in each domain. If no clustering is employed, then each predicate draws its value randomly from its domain, but when clustering is employed, the drawn values are clustered in certain regions. For example, when the number of clusters is 10, then all drawn values are centered (clustered) at 10 randomly chosen points in the domain. On the other hand, the cluster size models the extent of each individual cluster. For instance, the cluster size 400 implies that all values belonging to a particular cluster are within the Euclidean distance of 400. Studying these workload characteristics are essential when the data is being generated syntactically over a large domain.

k -index is highly sensitive to both the number of clusters and the cluster size. Both fewer clusters and smaller cluster size results in larger overlap among subscriptions and increases the

number of false candidates, especially for the Zipf distribution. Thus, from having no cluster, denoted by "N/A", to only two clusters, the matching time is increased by a factor of 20.84x, and from reducing the cluster size from 800 to 25, the matching time is also increased by a factor of 17.44x (Figure 7b,d).

Likewise, BE-Tree, by relying on a rigid clustering, is best suited when predicate values are uniformly drawn, namely, no data clustering, because the rigid clustering fails to fully adapt to the skewed distribution and fails to avoid indexing the dead spaces. As a result, BE*-Tree improves over BE-Tree by up to 77% when data is clustered and by up to 65% when no clustering is used (Figure 7b). Finally, BE*-Tree remains dominant as the number and size of clusters are varied, as observed in Figure 7.

Workload Size/Matching Probability (top-k) To enable top-k processing, we must first generate the predicate weights. In fact, our model is not limited to any particular weight generation technique or scoring function for top-k. Therefore, we adopt the weight generation technique proposed in [2]. In brief, for each unique attribute, we first compute its inverse frequency, denoted by $inv(attr)$, i.e., a popular attribute will be assigned a low weight while a rare attribute will be assigned a high weight. Second, for each expression Ω_i , we compute its predicate weight using a random Gaussian distribution (adopted from [2]) as follows:

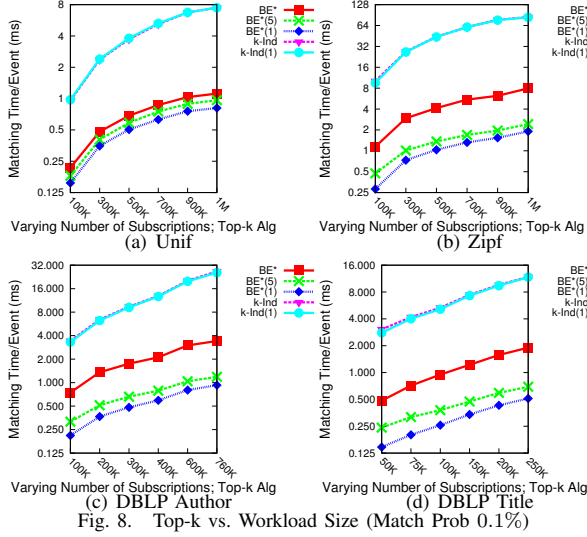


Fig. 8. Top-k vs. Workload Size (Match Prob 0.1%)

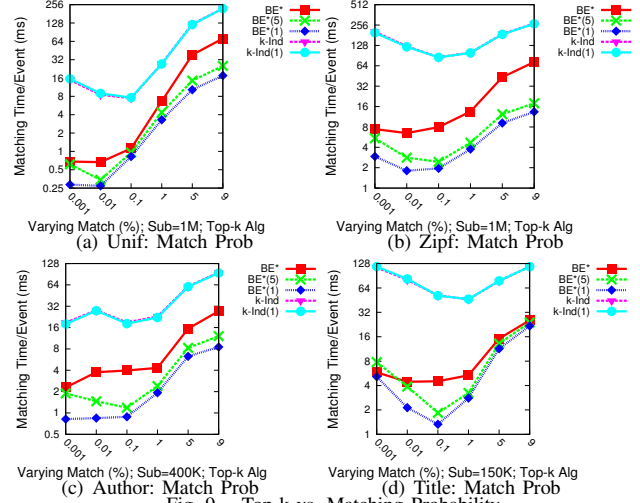


Fig. 9. Top-k vs. Matching Probability

$$\forall P_j(x) \in \Omega_i, P_j^{vt} = \max(\text{inv}(P_j^{\text{attr}}), \mathcal{N}(0.8 \times \text{inv}(P_j^{\text{attr}}), 0.05 \times \text{inv}(P_j^{\text{attr}})))$$

Next, we present BE*-Tree's top-k processing results for varying the workload size (Figure 8) and the matching probability (Figure 9). BE*-Tree hierarchical top-k processing exhibits a significant performance gain compared to k -index for finding the top-1 (denoted by $k\text{-ind}(1)$ and $\text{BE}^*(1)$) and top-5 ($\text{BE}^*(5)$) most relevant expression(s). BE*-Tree's success is attributed to its hierarchical structure that refines the upper bound score as traversing down the tree; however, k -index is restricted to a flat structure that is unable to further refine the upper bound score as more expressions are seen and limited to a fix top-layer upper bound scoring. In short, as we increase the workload size (Figure 8), $\text{BE}^*(1)$ improves over BE^* -Tree by up to 23%, 77%, 72%, 71% while $k\text{-ind}(1)$ improves over k -index by only up to 1.3%, 6.2%, 3.5%, 8.5%⁷ for the uniform, Zipf, author, and title datasets, respectively. Similar trends that show BE*-Tree's effectiveness are also observed as we increase the matching probability (Figure 9). In general, an increase in matching probability results in an increase in the number of candidate matches. Therefore, k -index is forced to scan large posting lists with a reduced chance of pruning and an increased application of sorting to advance through posting lists. The increased matching time is also observed in BE*-Tree simply due to an increased number of matches. However, BE*-Tree scales better as matching time increases and continues to outperform k -ind.

VIII. CONCLUSIONS

To address the problem of finding the most relevant matches, we present BE*-Tree, a data structure that employs an effective hierarchical top-k pattern matching algorithm. Furthermore, BE*-Tree introduces a novel non-rigid space-cutting technique to index Boolean expressions over a high-dimensional space by developing (1) a bi-directional tree expansion technique that

⁷Similar minor improvements were also reported in [2] that showed k -index's top-k achieved only up to 11% improvement over the baseline on some datasets while resulted in up to 19% performance loss on others.

enables indexing only non-empty continuous sub-spaces and (2) a descendant-repelling and overlap-free splitting strategy.

Moreover, BE*-Tree is a general index structure for efficiently processing large collections of Boolean expressions, a core functionality, required by a wide range of applications including (complex) event processing, computational finance, computational advertising, selective information dissemination, and location-based services. Finally, the performance of BE*-Tree is proven through a comprehensive experimental comparison with state-of-the-art index structures for matching Boolean expressions.

REFERENCES

- [1] A. Machanavajjhala, E. Vee, M. Garofalakis, and J. Shanmugasundaram, "Scalable ranked publish/subscribe," *VLDB*'08.
- [2] S. Whang, C. Brower, J. Shanmugasundaram, S. Vassilvitskii, E. Vee, R. Yemeni, and H. Garcia-Molina, "Indexing boolean expressions," in *VLDB*'09.
- [3] M. Fontoura, S. Sadanandan, J. Shanmugasundaram, S. Vassilvitskii, E. Vee, S. Venkatesan, and J. Zien, "Efficiently evaluating complex Boolean expressions," in *SIGMOD*'10.
- [4] M. K. Aguilera, R. E. Strom, D. C. Sturman, M. Astley, and T. D. Chandra, "Matching events in a content-based subscription system," in *PODC*'99.
- [5] F. Fabret, H. A. Jacobsen, F. Llirbat, J. Pereira, K. A. Ross, and D. Shasha, "Filtering algorithms and implementation for fast pub/sub systems," *SIGMOD*'01.
- [6] A. Campailla, S. Chaki, E. Clarke, S. Jha, and H. Veith, "Efficient filtering in publish-subscribe systems using binary decision diagrams," in *ICSE*'01.
- [7] M. Sadoghi and H.-A. Jacobsen, "BE-Tree: An index structure to efficiently match Boolean expressions over high-dimensional discrete space," in *SIGMOD*'11.
- [8] A. Guttman, "R-tree dynamic index structure for spatial searching," *SIGMOD*'84.
- [9] N. Beckmann, H.-P. Kriegel, R. Schneider, and B. Seeger, "The R*-tree: an efficient and robust access method for points and rectangles," *SIGMOD Rec.*'90.
- [10] T. Yan and H. Garcia-molina, "Index structures for selective dissemination of information under the boolean model," *ACM TODS*'94.
- [11] V. Gaede and O. Günther, "Multidimensional access methods," *ACM Comput. Surv.*'98.
- [12] R. Fagin, A. Lotem, and M. Naor, "Optimal aggregation algorithms for middleware," ser. PODS '01.
- [13] I. F. Ilyas, G. Beskales, and M. A. Soliman, "A survey of top-k query processing techniques in relational database systems," *ACM Comput. Surv.*'08.
- [14] M. Freeston, "A general solution of n-dimensional B-tree problem," *SIGMOD*'95.