

A Novel Similarity-Search method for Mathematical Content in *LaTeX* Markup

ABSTRACT

A relaxed structural matching search method, along with a symbolic similarity measurement algorithm for mathematical content search is proposed. Our approach uses an intermediate tree representation to capture structural information of mathematical expression, and based on a previous idea which indexes math expression structure through tree leaf-root paths, we further describe an advanced AND search method in a formal way. This search method can be used to test query/document subexpression isomorphism or evaluate the symbolic similarity between math expressions with consideration of their α -equivalence. For the purpose of evaluation, we also implement a search engine based on our idea.

Categories and Subject Descriptors

H.3 [Information Search and Retrieval]: Miscellaneous

General Terms

Algorithms

Keywords

mathematical searching, language processing, search engine

1. INTRODUCTION

With *MathJax* becoming popular, more and more \LaTeX markup can be crawled directly from many websites. In order to search those mathematical language in \LaTeX markup, a search method that can respect the properties of math expression needs to be developed. Although many researches have been conducted to retrieve information in structured content (e.g. *MathML*), information retrieval on \LaTeX math content is still not well-studied or exhaustively covered by mainstream IR research, compared to that on general text.

Unlike general text content, mathematical language, by its nature, has many differences from other textual document,

there are a number of new problems in measuring mathematical expression similarity. Among those problems, we know one math expression can be transformed to alternative forms, e.g. $\frac{a+b}{c}$ and $\frac{a}{c} + \frac{b}{c}$ should be considered as semantically identical. To identify those variations requires search engine to apply mathematical transformation rules to a query in order to obtain all forms of relevant expressions. Further, math expressions with the same evaluated value may also be considered relevant, in this case, $\sin(\frac{\pi}{2})$ and 1 are equivalent. Some computational search engines (e.g. *Symbolab* and *WolframAlpha*) are aware of these problems. But sometimes we need to rely on conventions and context to distinguish expressions such as $f(a+b)$ and $c(a+b)$, because the symbol f in the former expression is likely to represent function instead of a variable, in addition, expression such as f^{-1} can either be reciprocal or an inverse function. Moreover, a higher level of understanding of mathematic knowledge may be required for math-aware search engine to find the results for queries such as “find an article related to the four color theorem” [1].

Yet the problems addressed above are not considered in this paper, instead, we are focusing on the aspects which does not require a “good understanding” of mathematics, we target our research domain to be the following: The first is structural similarity. For example, $ax + (b + c)$ is not equivalent to $(a + b)x + c$ although they have the same set of symbols, this is because their structural difference. However, as the position of operands in math expression can be commutative in some cases, structural similarity is often measured by substructure isomorphism if we use operation tree [2] to represent math expressions. The second is symbolic similarity, with the consideration of α -equivalence. We know that symbols can be used interchangeably in each math formula to express the same meaning, e.g. $a^2 + b^2 = c^2$ and $x^2 + y^2 = z^2$. Nevertheless, we still weight symbolic similarity sometimes, for instance, $E = mc^2$ is considered more meaningful when exact symbols are used rather than just being structurally identical with $y = ax^2$. On the other hand, we should also weight α -equivalent expressions more, that is, changes of symbols in expression preserve more syntactic similarity when changes are made by substitution, e.g. for query $x(1 + x)$, expression $a(1 + a)$ are considered more relevant than $a(1 + b)$. Because the “bond variable” x and a here are at the same positions and both supposed to represent the same value. All the points addressed here makes transitional IR methods (e.g. bag of words model and tf-idf weighting) deficient to handle math content.

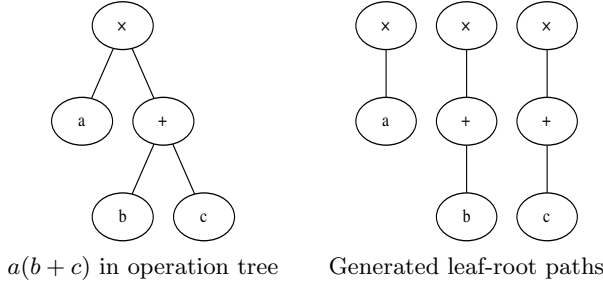


Figure 1: Leaf-root path example

2. RELATED WORK

Similarity/boolean search for mathematical content is not a new topic, conference in this topic is getting increasingly research attention and the proposed systems have progressed considerably [3]. DLMF project from NIST [4, 5] and MIA system [6, 7, 8], notably, use text-based approaches and utilize existing models to deal with math content on top of existing IR tools (such as *Apache Lucene*). They are commonly using augmentation and normalization (by ordering the subexpressions) to enumerate and represent all possible sequences of commutative operands (symbols (e.g. a and b) and unified items (id , $const$) in different levels of math expression. MWS [9, 10, 11] takes an *automatic theorem proving* approach and uses *term indexing* to minimize the cost of unification algorithm which is able to find if two expressions are equivalent, however, their index relies on RAM memory [11] and needs to include all sub-terms of a formula [9]. A *symbol layout tree* or *presentation tree* [12, 13] is introduced to describe geometric layouts of symbols in a formula. [13] uses two templates to parse L^AT_EX markup with two typical operator terms: explicit ones (“ $\frac{}{}$ ”, “ $\sqrt{}$ ”, etc.) and implicit ones (“+”, “ \div ”, etc.) to form a presentation tree, then extracts original terms and generalized terms from normalized presentation tree, to provide the flexibility of both fuzzy and exact search. [12] uses symbol layout tree as a kind of substitution tree, while [14, 15] have developed a *symbol pairs* idea to capture relative position information between symbol pairs, which enables key-value storage to speed search. [16] extracts HASH signature for each subtree, their system has a slower index growth resulted from a lot of occurrences of common subtree structure in math expressions. [17] extracts semantic keywords, structure description to indicate subordinative relationship in a string format. The similarity is calculated using normalized tf-idf vector. Tree edit distance is adopted by [18, 19] in which they try to overcome the bad time complexity of original algorithm by summarizing and using a compromised edit distance algorithm, then by applying with an early termination algorithm along with a distance cache. There are also efforts using image-based approaches [20, 21] and lattice-based approach [22] to measure math formula similarity in a different perspective.

3. METHODOLOGY

Our method can be seen as an approach built upon the idea of leaf-root path or sub-path [23, 24, 25, 26] from an operation tree, to capture structure and semantic information of math expression. Figure 1 is an example of generated leaf-root paths for math expression $a(b+c)$. The intuition behind

this idea is that an operation tree, no matter how operands are ordered, uniquely determines the leaf-node paths decomposed from the tree. This makes leaf-root path a good fit for representing mathematical expression because commutative operands are exhaustively used in mathematical language. Besides, by going bottom-up from leaves of an operation tree, we are essentially traversing to an expression from its subexpression for every level. So we can index the leaf-root paths and search an expression by going through and beyond the leaf-root paths from its subexpression.

We develop these ideas to simultaneously search along the way of all leaf-root paths, so that we are essentially pruning indexes which does not share the common postfixes beyond the root of the query tree. To better describe further ideas build upon this, we will put it in a formal way.

3.1 Formal Definition

Here we clarify some notations used throughout this paper, a path p is a sequence of numbers given by $p = p_0 p_1 \dots p_n$, $n \geq 0$, $p_i \in \mathbf{R}$ and \mathbf{P} is the set of all paths. Any function $f(\cdot) = y \in \mathbf{R}$ applied on path p is mapped to a path too: $f(p) = f(p_0)f(p_1)\dots f(p_n)$. And we name a *concatenation* of two paths $^1p = p_0 p_1 \dots p_n$ and $^2p = p_n p_{n+1} \dots p_m$ where $m \geq n$, to be a new path $^1p \cdot ^2p$ defined as $p_0 p_1 \dots p_n p_{n+1} \dots p_m$, and the concatenation of a path p on a set $S = \{s_1, s_2 \dots s_n\}$ is defined as $S \cdot p = \{s_1 \cdot p, s_2 \cdot p \dots s_n \cdot p\}$. A path with only one element can be explicitly wrapped by a bracket, e.g. $p = (p_0)$, to avoid confusion. Furthermore, the *longest common postfix* path p^* between two path p_1 and p_2 is mapped by the function named lcp, which is defined by $p^* = \text{lcp}(p_1, p_2) = \text{lcp}(p_2, p_1)$.

We introduce a *formula tree* to represent a mathematical expression, in which each node is associated with a label to represent the unified token (e.g. same value for token +, \oplus and \pm) and each leaf node is associated with a symbol to identify math expression operand. Besides, a *formula subtree* relation is also defined to address the sub-structure relation between two mathematical expressions.

3.1.1 Formula tree

A *formula tree* is a labeled rooted tree $T = T(V, E, r)$ with root r , where each vertices $v \in V(T)$ is associated with a label (not necessarily unique in the same tree) $\ell_T(v) \in \mathbf{R}$ mapped by label function ℓ_T , and each leaf $l \in V(T)$ is also associated with a symbol $S_T(v) \in \mathbf{R}$ mapped by symbol function S_T . For convenience, we will write ℓ and S as short names which refer to the tree implied by the context, also we use function $S(p)$ to indicate the symbol of the leaf in a leaf-root path p .

3.1.2 Formula subtree

Given formula tree S and T , we say S is a *formula subtree* of T if there exists an injective mapping $\phi : V(S) \rightarrow V(T)$ satisfying:

1. $\forall (v_1, v_2) \in E(S)$, we have $(\phi(v_1), \phi(v_2)) \in E(T)$;
2. $\forall v \in V(S)$, we have $\ell(v) = \ell(\phi(v))$;
3. If $v \in V(S)$ is a leaf vertices in S , then $\phi(v)$ is also a leaf in T .

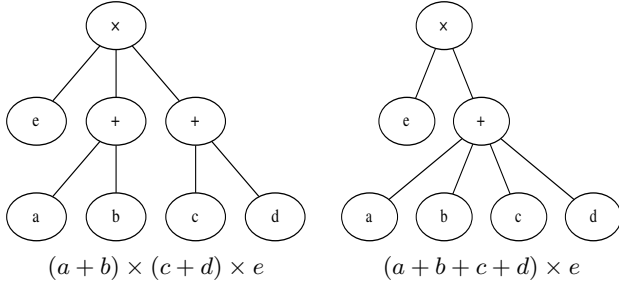


Figure 2: Leaf-root paths with different structure

Such a mapping ϕ is called a formula subtree isomorphic embedding (or formula embedding) for $S \rightarrow T$. If satisfied, we denote $S \preceq_l T$ on Φ , where Φ ($\Phi \neq \emptyset$) is the set of all the possible formula embeddings for $S \rightarrow T$.

3.1.3 Leaf-root path set

A *leaf-root path set* generated by tree T is a set of all the leaf-root paths from tree T , mapped by a function $g(T)$.

3.1.4 Index

An *index* Π is a set of trees such that $\forall T \in \Pi$, we have $T \in \mathcal{I}_\Pi(a)$ for any $a \in \ell(g(T))$, we say T is indexed in Π and \mathcal{I}_Π is called an index look-up function for index Π .

3.2 Search Method

For a collection of document expressions, we will index them by merging all the leaf-root paths from each document formula tree into a large “inverted” index tree, in which each node from path a stores the information of all the indexed formula trees in $\mathcal{I}_\Pi(a)$.

Through searching all sub-paths at the same time, we are able to limit the set of formula trees being structurally matching (in formula subtree relation) with a query formula tree, to only a subset of our index. This is illustrated as follows.

Given an index Π and a formula tree T_q , $\forall T_d \in \Pi$: If $T_q \preceq_l T_d$ on Φ , then $\exists \hat{a} \in \mathbf{P}$, s.t.

$$T_d \in \bigcap_{a \in L} \mathcal{I}_\Pi(a)$$

where $L = \ell(g(T_q)) \cdot \hat{a}$.

Justification. Denote the root of T_q and T_d as r and s respectively. Let \hat{p} be the path determined by vertices from $t = \phi(r)$ to s in T_d , and ${}^1p, {}^2p \dots {}^np$, $n \geq 1$ be all the leaf-node paths in T_q . Then $\hat{a} = \ell(\hat{p})$, this is because: $L = \ell(\{{}^1p, {}^2p \dots {}^np\}) \cdot \hat{a} = \ell(\{\phi({}^1p), \phi({}^2p) \dots \phi({}^np)\}) \cdot \ell(\hat{p}) = \{\ell(\phi({}^1p) \cdot \hat{p}), \ell(\phi({}^2p) \cdot \hat{p}) \dots \ell(\phi({}^np) \cdot \hat{p})\}$. According to definition 3.1.2 and $t = \phi(r)$, we have $\phi({}^ip) \cdot \hat{p} \in g(T_d)$, $1 \leq i \leq n$. Since $T_d \in \Pi$, T_d is indexed in Π with respect to each of the elements in L , that is to say $\forall a \in L$, $T_d \in \mathcal{I}_\Pi(a)$.

In a nutshell, we search the index by intersecting the indexed formula trees from all the generated leaf-node paths at the same time, then further possible search path \hat{a} is only possible when paths along the generated leaf-node paths in the index have a common postfix. Therefore we can “merge” the

```

1: procedure REMOVECANDIDATE( $d, Q, C$ )
2:   for  $a \in Q$  do
3:     if  $C_a = \emptyset$  then
4:       return  $\emptyset$ 
5:     else
6:        $C_a := C_a - \{d\}$ 
7:   return  $C$ 
8:
9: procedure MATCH( $a, a', Q, C$ )
10:  for  $b \in Q$  do
11:     $t := \text{lcp}(a, b)$ 
12:     $Q_t := Q_t \cup \{b\}$ 
13:     $P := P \cup \{t\}$ 
14:  for  $t \in P$  do
15:    for  $b \in Q_t$  do
16:      for  $b' \in C_b$  do
17:        if  $t \neq \text{lcp}(a', b')$  then
18:           $C := \text{REMOVECANDIDATE}(b', Q_t, C)$ 
19:        if  $|C| = 0$  then
20:          return FAIL
21:    if  $\text{DECOMPOSEANDMATCH}(Q_t, C) = \text{FAIL}$  then
22:      return FAIL
23:  return SUCC
24:
25: procedure DECOMPOSEANDMATCH( $Q, C$ )
26:  if  $Q = \emptyset$  then return SUCC
27:   $a := \text{OnePathIn}(Q) \triangleright$  Choose a reference path in  $Q$ 
28:   $Q_{\text{new}} := Q - \{a\}$ 
29:  for  $a' \in C_a$  do
30:     $C_{\text{new}} := \text{REMOVECANDIDATE}(a', Q_{\text{new}}, C)$ 
31:    if  $C_{\text{new}} = \emptyset$  then return FAIL
32:    if  $\text{MATCH}(a, a', Q_{\text{new}}, C_{\text{new}})$  then return SUCC
33:  return FAIL

```

Figure 3: The decompose-and-match algorithm

paths ahead and prune those paths not in common. Level by level, we are always able to find the structurally matched formula tree as long as it is indexed in Π .

3.3 Substructure Matching

However, query formula tree will not necessarily being formula subtree to all the document (indexed) formula trees in our search set $\bigcap_{a \in L} \mathcal{I}_\Pi(a)$, even if their generated leaf-root paths are identical. One supporting example for this point is shown in figure 7. To address this problem, we propose an algorithm described in in figure 3, to test the document formula trees in our search set to see if they are in formula subtree relation. This algorithm is inspired from the following observations.

3.3.1 Observation 1

For two formula trees which satisfy $T_q \preceq_l T_d$ on Φ , then $\forall \phi \in \Phi$, $p \in g(T_q)$, also any vertices v along path p , the following properties are obtained:

$$\deg(v) \leq \deg(\phi(v)) \quad (1)$$

$$\ell(p) = \ell(\phi(p)) \quad (2)$$

$$|g(T_q)| \leq |g(T_d)| \quad (3)$$

Justification. Because $\forall w \in V(T_q)$ s.t. $(v, w) \in E(T_q)$, there exists $(\phi(v), \phi(w)) \in E(T_d)$. And for any (if exists) two

different edges $(v, w_1), (v, w_2) \in E(T_q)$, $w_1 \neq w_2 \in V(T_q)$, we know $(\phi(v), \phi(w_1)) \neq (\phi(v), \phi(w_2))$ by definition 3.1.2. Therefore any different edge from v is associated with a distinct edge from $\phi(v)$, thus we can get (1). Given the fact that every non-empty path p can be decomposed into a series of edges $(p_0, p_1), (p_1, p_2) \dots (p_{n-1}, p_n)$, $n > 0$, property (2) is trivial. Because there is exact one path between every two nodes in a tree, the leaf-root path is uniquely determined by a leaf node in a tree. Hence the rationale of (3) can be obtained in a similar manner with that of (1), expect neighbor edges are replaced by leaf-node paths.

3.3.2 Observation 2

Given two formula trees T_q and T_d , if $|g(T_q)| = 1$ and $\ell(g(T_q)) \subseteq \ell(g(T_d))$, then $T_q \preceq_l T_d$.

Justification. Obviously there is only single one leaf-root path in T_q because $|g(T_q)| = 1$. Denote the path as $p = p_0 \dots p_n$, $n \geq 0$ where p_n is the leaf, and let $a = \ell(p)$. Since $a \subseteq \ell(g(T_d))$, we know that there must exist a path $p' = p'_0 \dots p'_n \in g(T_d)$ such that $a = \ell(p')$. Without loss of generality, suppose p'_n is the leaf of T_d . Now the injective function $\phi : p_i \rightarrow p'_i$, $0 \leq i \leq n$ satisfies all the requirements for T_q as a formula subtree of T_d .

3.3.3 Observation 3

For two formula trees T_q and T_d , if $T_q = T(V, E, r) \preceq_l T_d$ on Φ , $\forall a, b \in g(T_q)$ and a mapping $\phi \in \Phi$. Let $T'_d = {}^t T_d$ where $t = \phi(r)$ and $a' = \phi(a)$, $\forall b' \in g(T'_d)$, it follows that:

$$b' = \phi(b) \Rightarrow |\text{lcp}(a, b)| = |\text{lcp}(a', b')|$$

Furthermore, $\forall c \in g(T_q)$ s.t. $|\text{lcp}(a, b)| \neq |\text{lcp}(a, c)|$, we have

$$|\text{lcp}(a, b)| = |\text{lcp}(a', b')| \Rightarrow b' \neq \phi(c)$$

Justification. Because $a, b \in g(T_q)$, thus $a_0 = b_0 = r$, and we can also make sure $\text{lcp}(a, b) \geq 1$. Denote the path of $a = a_0 \dots a_n a_{n+1} \dots a_{l-1}$, similarly denote the path of b as $b = b_0 \dots b_m b_{m+1} \dots b_{m-1}$, where the length of each $l, m \geq 1$ and $a_i = b_i$, $0 \leq i \leq n \leq \min(l-1, m-1)$ while $a_{n+1} \neq b_{n+1}$ if $l, m > 1$. On the other hand $a' = \phi(a)$ and $b' \in g({}^t T_d)$, therefore $a'_0 = \phi(a_0) = \phi(r) = t = b'_0$. For the first conclusion, if $b' = \phi(b)$, there are two cases. If any of $|a|$ or $|b|$ is equal to one then $|\text{lcp}(a, b)| = |r| = |t| = |\text{lcp}(a', b')| = 1$; Otherwise if $l, m > 1$, path $a_0 \dots a_n = b_0 \dots b_n$ and $a_{n+1} \neq b_{n+1}$ follow that $\phi(a_0 \dots a_n) = \phi(b_0 \dots b_n)$ and $\phi(a_{n+1}) \neq \phi(b_{n+1})$ by definition. Because edge $(\phi(a_n), \phi(a_{n+1}))$ and $(\phi(b_n), \phi(b_{n+1}))$ are also in $E(T'_d)$, we have $|\text{lcp}(a, b)| = |\text{lcp}(a', b')| = n$. For the second conclusion, we prove by contradiction. Assume $b' = \phi(c)$, by the first conclusion we know $|\text{lcp}(a, c)| = |\text{lcp}(a', b')|$. On the other hand, because $|\text{lcp}(a, c)| \neq |\text{lcp}(a, b)| = |\text{lcp}(a', b')|$, thus $|\text{lcp}(a, c)| \neq |\text{lcp}(a', b')|$ which is impossible.

For a query formula tree T_q and a document formula tree T_d , observation 1 offers us some constrains for finding initial possible isomorphic paths in T_d (what we call *candidates*) for a given query path in T_q . And observation 2 is a sufficient condition to test substructure relation, but the query tree has to consist only one leaf-root path to be tested. While observation 3 states two necessary conditions for one formula tree to be a formula subtree of another and implies that a group of query leaf-root paths can only be isomorphic to

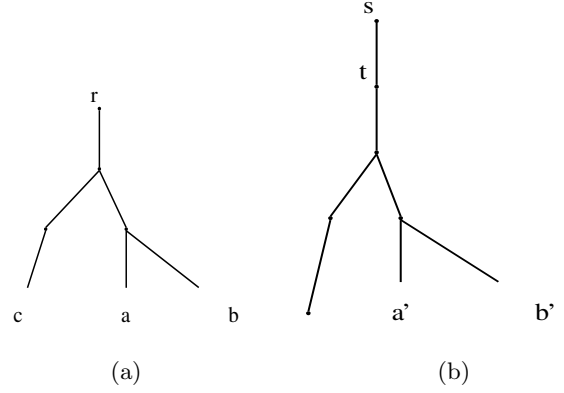


Figure 4: Formula subtree matching

someones in another group of leaf-root paths in document formula tree. This leads to the idea to decompose the formula tree and divide the problem into subproblems by ruling out impossible candidates using observation 3, until at some point we can apply observation 2 or observe trivial cases to solve all the sub-matching problems.

To illustrate our DECOMPOSEANDMATCH algorithm, figure 4 gives a general case in which query tree in (a) is trying to match a document tree in (b). Initially every leaf-root path in (a) should be associated with a set of candidate paths in (b) which satisfy the constrains in 3.3.1. Then we arbitrarily choose a path in (a) as a *reference path* (heuristically a *heavy path* [27]), for each of the paths in its candidate set, we choose it as reference path in (b), and suppose we choose a' here. At this time we can apply the two constrains from 3.3.3 and ruling out some impossible isomorphic paths in candidate set of each path in (a) and divide the problems further. For example, because $|\text{lcp}(a, b)| = |\text{lcp}(a', b')|$, we know b' is still in candidate set of b ; while b' is not in candidate of c because $|\text{lcp}(a, b)| \neq |\text{lcp}(a, c)|$. After going through these eliminations for each leaf-node path (except the reference path a) in (a), we now have two similar subproblems: c as a subtree along with its candidate set, and b as a subtree along with its candidate set. We can apply this algorithm recursively until a trivial subproblem is reached (e.g. the case in 3.3.2). During this process, if we find any candidate set to be empty, we stop the subproblem process and change to another reference path or stop the algorithm completely if every possible reference path is tried. The argument Q and C is the set of leaf-root paths in query tree and the candidate sets associated with all leaf-root paths respectively. The procedure returns SUCC if a matching possibility is found, otherwise FAIL is returned indicating the formula tree in (a) cannot not be a formula subtree of that in (b).

DECOMPOSEANDMATCH algorithm offers a way to “double-check” structure isomorphism, because for any two formula tree T_q and T_d and $\forall \hat{a} \in \mathbf{P}$, if $\ell(g(T_q)) \cdot \hat{a} \subseteq \ell(g(T_d))$, it is not sufficient to imply $T_q \preceq_l T_d$. Nevertheless, we think the cases which makes the above statement insufficient are fairly rare in common mathematical content, and the complexity introduced from this algorithm will offset the benefit to identify the structure isomorphism. Thus a compromised search,

for efficiency reason, would assume all the document formula trees in search set $\bigcap_{a \in L} \mathcal{I}_\Pi(a)$ is structurally matching the given query formula tree T_q .

3.4 Structural Similarity

After searching structurally matching document expressions in a boolean manner, we use two factors to measure their structural similarity degree.

The first is *matching depth*. As it is addressed in [6], the deeper sub-formulae in in mathematical expression will make it less important to the overall formula. For example, given query formula \sqrt{a} , expression \sqrt{x} would score higher than $\sqrt{\sqrt{x}}$ does. To reflect the depth where two expressions match, we define *matching depth factor* $f(d)$ to be a function value in negative correlation with matching depth $d = |\hat{a}|$, e.g. $f(d) = 1/(1+d)$ in our method. And use this factor in our final ranking process. The second is *matching ratio*. According to the property (2) in 3.3.1, $r = |g(T_q)|/|g(T_d)| \leq 1$, and the ratio r on left-hand is defined as *matching-ratio*, which characterises the structural coverage for the matching part in an expression. Intuition behind this is, for example, given query $\alpha y + \beta$, document expression $ax + b$ should precede $x^2 + ax + b$ simply because the query matches more “area” of the former document than that of the latter.

3.5 Symbolic Similarity

As we have discussed in section 1, besides structural similarity, symbolic similarity is also essential to be considered. Here our scoring goal for symbolic similarity can be summarized as:

- Given two formula trees $T_q \preceq_l T_d$ on Φ , if path $a \in g(T_q)$ is isomorphic to path $a' \in g(T_d)$, that is to say, $\phi(a) = a'$, where $\phi \in \Phi$, then if their symbol matches, i.e. $\mathcal{S}(a) = \mathcal{S}(a')$, we score them higher than those do not match symbolically. And the more symbolic matches there are, the higher symbolic relevance degree it has.
- α -equivalent expressions have more symbolic relevance degree than those are not, and the more bond variables (at the structurally matching positions) two expressions match, the more symbolically relevant they are considered to be.

The two are illustrated as follows. Let the rank of a structural relevant math expression d be $r(d)$, and given query $\sqrt{a}(a-b)$ for instance. Then the first goal is essentially saying:

$$r(\sqrt{a}(a-b)) > r(\sqrt{a}(a-x)) > r(\sqrt{x}(x-y))$$

The second one has two symbols matching while the third one has no symbolic match at all.

By the second goal, we know:

$$r(\sqrt{x}(x-b)) > r(\sqrt{x}(y-b))$$

Because the second one does not have the same number of bond-variable matches as the first one does.

```

1: procedure MARKANDCROSS( $D, Q, C$ )
2:   score := 0
3:   if  $D = \emptyset$  then
4:     return 0
5:   for  $a' \in D$  do
6:      $T_{a'} := \text{unmark}$ 
7:   for  $v \in \mathcal{V}(D)$  do
8:      $B_v := 0$ 
9:   QList := SORTBYSYMBOLANDOCUR(Q)
10:  for  $a$  in QList do
11:    for  $v \in \mathcal{V}(D)$  do
12:       $m := -\infty$ 
13:       $m_p := \emptyset$ 
14:      for  $a' \in C_a \cap \{y \mid y = v, y \in \mathcal{V}(D)\}$  do
15:        if  $T_{a'} = \text{unmark}$  and  $\text{sim}(a, a') > m$  then
16:           $m := \text{sim}(a, a')$ 
17:           $m_p := a'$ 
18:        if  $m_p \neq \emptyset$  then
19:           $T_{m_p} := \text{mark}$ 
20:           $B_v := B_v + m$ 
21:        else ▷ Exhausted all candidates
22:          return 0
23:      if  $\mathcal{S}(a)$  changed or last iteration of  $a$  then
24:         $m := -\infty$ 
25:         $m_v := \emptyset$ 
26:        for  $v \in \mathcal{V}(D)$  do
27:          if  $B_v > m$  then
28:             $m := B_v$ 
29:             $m_v := v$ 
30:             $B_v := 0$ 
31:        score := score +  $m$ 
32:        for  $v \in \mathcal{V}(D)$  do
33:          if  $v = m_v$  then
34:            nextState := unmark
35:          else
36:            nextState := cross
37:          for  $a' \in C_a \cap \{y \mid y = v, y \in \mathcal{V}(D)\}$  do
38:            if  $T_{a'} = \text{mark}$  then
39:               $T_{a'} := \text{nextState}$ 
40:  return score

```

Figure 5: The mark-and-cross algorithm

However, sometimes the two goals can be conflicting. Given document expression $\sqrt{a}(x-b)$ and $\sqrt{x}(x-b)$ for instance, the former has two symbolic matches (i.e. “ a ” and “ b ”) while does not have bond-variable match. The latter, on the other hand, has bond-variable match while only has one symbolic match (i.e. “ b ”). We nevertheless score the latter higher because it does not lose any mathematic semantics.

To meet the goals above, intuitively, we first take the bond variable with greatest number of occurrence in its expression in the query, try to match as many symbols as possible with bond variables from document expression. The bond variable with most matches would be chosen (named *best-matching*) to contribute to final symbolic relevance score (proportionally to the number of matches in that bond variable), and we exclude its path from matching query paths in future iterations. In the next iteration, we choose the bond variable with the second number of occurrence and repeat this process until all the query bond variable are iterated.

We describe our algorithm in figure 5 which measures symbolic similarity given two expressions. The MARKANDCROSS algorithm takes three arguments, the set of leaf-root paths D and Q in document expression and query expression respectively, and the candidate sets C associated with all leaf-root paths in query. The bond variables in D is addressed by the set $\mathcal{V}(D) = \{x \mid \mathcal{S}(x), x \in D\}$, which contains all the leaf node symbols from document expression. Function SORTBYSYMBOLANDOCUR takes the elements from a set of leaf-root paths and returns a list containing all the paths. The list is sorted by the tuple $(\mathcal{S}(p), O_p)$, where O_p is the number of occurrence of path p in the list. Each document path a' is associated with a tag $T_{a'}$ which has three possible states: marked, unmarked and crossed. And bond variable $v \in \mathcal{V}(D)$ can be given a score B_v which represents the similarity degree between current evaluating query/document bond variables. The symbolic similarity function $\text{sim}(a, a')$ measures the similarity degree between two leaf-root paths a and a' . Intuitively, we set the similarity function:

$$\text{sim}(a, a') = \begin{cases} 1 & \text{if } \mathcal{S}(a) = \mathcal{S}(a') \\ \alpha < 1 & \text{otherwise} \end{cases}$$

to give more weights to leaf-root paths with exact symbol match.

Let us determine the proper value for α . Consider the conflicting cases stated in this section by using another example here, given query expression $a + \frac{1}{a} + \sqrt{a}$ and document expression $a + \frac{1}{a} + b + \frac{1}{b} + \sqrt{b}$, we consider bound variable matching

$$\boxed{a} + \frac{1}{\boxed{a}} + \sqrt{\boxed{a}}$$

with

$$a + \frac{1}{a} + \boxed{b} + \frac{1}{\boxed{b}} + \sqrt{\boxed{b}}$$

weighted more than exact symbol matching

$$\boxed{a} + \frac{1}{\boxed{a}} + \sqrt{a}$$

with

$$\boxed{a} + \frac{1}{\boxed{a}} + b + \frac{1}{b} + \sqrt{b}$$

(expressions surrounded by a box here indicates the matching part)

Because the former matching has more variables involved even if they are not identical symbolic matches compared with its counterpart of the letter. That is to say, given a bound variable matching k variables with that bound variable in query, we need α to satisfy $k\alpha > (k-1) \times 1 = k-1$ and $\alpha < 1$. Therefore, we set α to a value close to 1 (e.g. 0.9 in our practice).

By sorting the query paths Q , the algorithm is able to take out paths from same bound variable in maximum-occurrence-first order from QList. Each query path a tries to match a path a' in each document bond variable v by selecting the unmarked path m_p with maximum $\text{sim}(a, a')$ value, and accumulate the value on B_v indicating the similarity between

currently evaluating query bond variable and the bond variable v . In addition, mark the tag T_{m_p} associated with the document path m_p . Once a query bond variable has been iterated completely (line 23), we find the document variable m_v with greatest B_v value m , and regard it as the best-matching bond variable in $\mathcal{V}(D)$ for the bond variable just iterated, and add m to the result score. Before iterating a new query bond variable, we will cross all the document paths of variable m_v to indicate they are confirmed been matched, and unmark the tags of those marked paths that are not variable m_v . We continue doing so until all the query path is iterated, and finally return the result score indicating the symbolic similarity between the two expression.

4. IMPLEMENTATION

In order to evaluate the performance of our method, we have implemented a proof-of-concept search engine¹ as well as a parser for parsing L^AT_EX markup content directly into operation trees.

4.1 Parser

We are tokenizing math content using lexer generator *flex* and have implemented a LALR parser generated from a set of *GNU bison* grammar rules we defined specifically for MathJax content in a subset of L^AT_EX (those related to mathematics). Our parser transforms a math formula into an in-memory operation tree (representing formula tree), as an intermediate step to extract the path labels, path ID, and degree numbers associated with the path for every leaf-root path from the tree. Our lexer omits all L^AT_EX control sequences not matching any pattern of our defined tokens, most of them are considered unrelated to math formula semantics (environment statement, color, mbox etc.).

4.2 Index

The indexer writes the information extracted from parser into disk. There are two parts in our index, the first part uses native file system (for the sake of implementation simplicity) to store leaf-root path labels in directories from which our search engine can go level by level. Path ID, degree numbers, and also the formula ID generating that leaf-root path (we refer to these three as *branch word*) are stored in a “posting” file at the directory corresponding to that branch word labels, e.g. the example in figure 1 will index into two directories: `./VAR/TIMES` and `./VAR/ADD/TIMES`. All the indexed branch word with path labels corresponding to directory `./VAR/TIMES/` are stored in the posting file of that directory, locating at: `./VAR/TIMES/posting.bin`. Branch words in a posting file are ordered by their formula IDs to speed merge search. The second part of our index is a key-value database to map a formula ID to additional information for that formula (e.g. original markup, number of leaf-root paths $|T_d|$ and the URL on which the formula is crawled).

4.3 Searching and Ranking

The compromised substructure searching in indexed expressions (described in section 3.2) is used to filter out likely isomorphic expressions in our index. Searching is performed by simultaneously going from all the directories corresponding to the generated leaf-root paths from query formula tree,

¹demo page: <http://infolab.ece.udel.edu:8912/cowpie/>

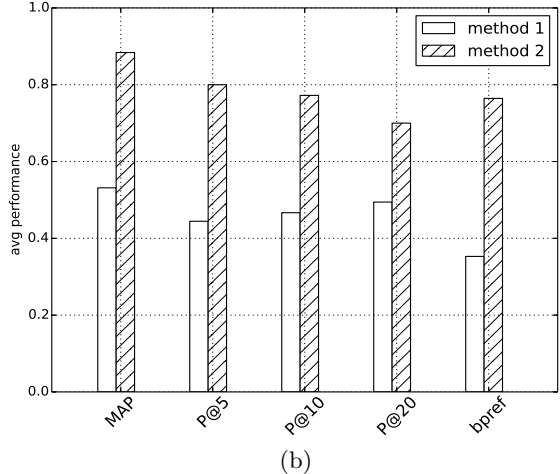
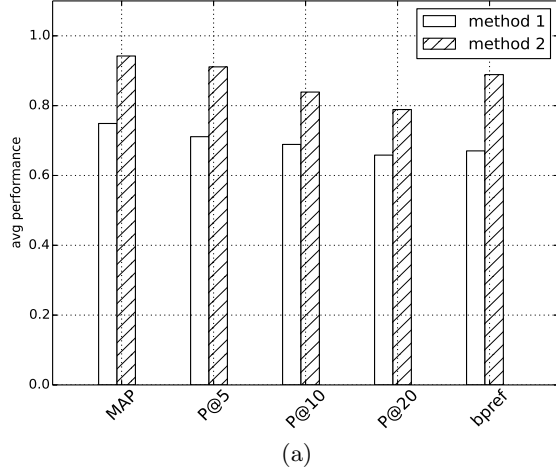


Figure 6: Effectiveness Comparison

to all their merged subdirectories. We go on traversing and intersect the branch words in the posting files of each searching directory by their formula IDs. All the intersected formula IDs actually represent the indexed formula trees in our search set $\bigcap_{a \in L} \mathcal{I}_{\Pi}(a)$. Every document formula identified in the search set is considered as a hit, we then apply MARKANDCROSS algorithm to get its symbolic similarity score with query formula. Denote this score to be s , and $f(d)$ being the matching depth factor, r being the matching ratio for one pair of query/document. We will use them together to indicate the overall similarity and rank items in search results, measured by tuple $(s, f(d), r)$. To know if one should ranked higher than another, first compare s , if equal, compare $f(d)$ and then r . In addition, we use a min-heap to keep the top- k scored items in our search results (by replacing the lowest scored items if we find a newer hit with higher score), where k is the maximum number of items a search results may contain. We also place a valve on the number of branch words can be searched for one query at a time, exceeding this limit will result in a stop of searching and immediately returning the search results we have so far. Because some query can potentially have a very long posting list, doing so would make our search response no more than a certain time.

5. EVALUATION

The popular evaluation dataset in this research domain, the NTCIR Math Task dataset, is in MathML/XML format, and original L^AT_EX information is not always preserved in their dataset. Because we are parsing L^AT_EX directly, although converting MathML/XML formula back into L^AT_EX is possible, we fail to convert all the document correctly (using *pandoc*). Furthermore, supporting wildcard query is a default requirement in NTCIR Math Task, while our approach does not find a way to enable wildcard feature so far. Nevertheless, we choose to create our own dataset to evaluate our method. We have crawled pure L^AT_EX content from the posts of nearly entire Math Stack Exchange website

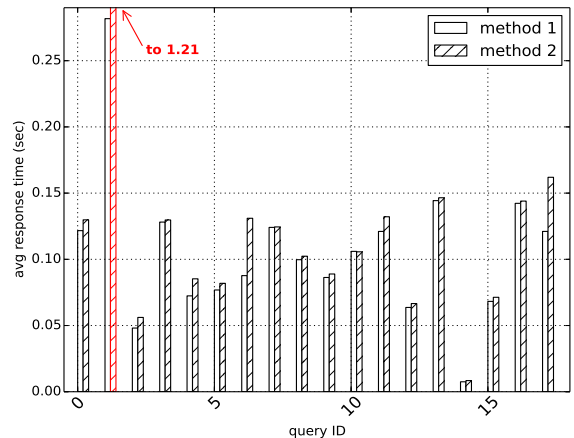


Figure 7: Efficiency Performance

before March 2015 (over 60MB *bzip2* compressed raw data² in plain text format with one L^AT_EX math mode content per line, and one file for each post). Our test query set³ consists queries mostly from [1] and [14], some of them are excluded here because we are not able to find similar formula in our own dataset. Table 1 shows our complete test queries used in our evaluation.

There are four levels of relevance scored from 0 to 4 in our evaluation. The criteria considers both structural similarity and symbolic similarity. Structural similarity is measured by 0, 1 (mostly similar) and 2 (complete matching), and symbolic similarity is measured by 0, 1 (mostly identical symbols for the matching parts) and 2 (identical symbols for the matching parts). The level of relevance is simply the sum

²<http://infolab.ece.udel.edu:8912/cowpie/math-stack-exchange-march-2015.tar.bz2>

³<http://infolab.ece.udel.edu:8912/cowpie/test-queries.tex>

of the two scores. We have evaluated our method (method 2) and a baseline method (method 1) without considering symbolic similarity. Table 2 shows the distribution of hits and relevance level for each run.

6. CONCLUSION

There are four levels of relevance scored from 0 to 4 in our evaluation. The criteria considers both structural similarity and symbolic similarity. Structural similarity is measured by 0, 1 (mostly similar) and 2 (complete matching), and symbolic similarity is measured by 0, 1 (mostly identical symbols for the matching parts) and 2 (identical symbols for the matching parts).

7. REFERENCES

- [1] Topics for the ntcir-10 math task full-text search queries. <http://ntcir-math.nii.ac.jp/wp-content/blogs.dir/13/files/2014/02/NTCIR10-math-topics.pdf>.
- [2] Richard Zanibbi and Dorothea Blostein. Recognition and retrieval of mathematical expressions. *International Journal on Document Analysis and Recognition (IJ DAR)*, 15(4):331–357, 2012.
- [3] Akiko Aizawa, Michael Kohlhase, and Iadh Ounis. Ntcir-11 math-2 task overview. *The 11th NTCIR Conference*, 2014.
- [4] Miller B. and Youssef A. Technical aspects of the digital library of mathematical functions. *Annals of Mathematics and Artificial Intelligence* 38(1-3), 121-136, 2003.
- [5] Youssef A. Information search and retrieval of mathematical contents: Issues and methods. *The ISCA 14th IntâĂłl Conf. on Intelligent and Adaptive Systems and Software Engineering (IASSE 2005)*, 2005.
- [6] Petr Sojka and Martin L. Indexing and searching mathematics in digital libraries. *Intelligent Computer Mathematics*, 6824:228–243, 2011.
- [7] Petr Sojka and Martin L. The art of mathematics retrieval. *ACM Conference on Document Engineering, DocEng 2011*, 2011.
- [8] Martin L. Evaluation of mathematics retrieval. Master’s thesis, Masarykova University, 2013.
- [9] Michael Kohlhase and Ioan A. SăĂşucan. A search engine for mathematical formulae. In *Proc. of Artificial Intelligence and Symbolic Computation, number 4120 in LNAI*, pages 241–253. Springer, 2006.
- [10] Michael Kohlhase. Mathwebsearch 0.4 a semantic search engine for mathematics.
- [11] Michael Kohlhase. Mathwebsearch 0.5: Scaling an open formula search engine.
- [12] Thomas Schellenberg, Bo Yuan, and Richard Zanibbi. Layout-based substitution tree indexing and retrieval for mathematical expressions. *Proc. SPIE 8297, Document Recognition and Retrieval XIX*, 82970I, 2012.
- [13] Xuan Hu, Liangcai Gao, Xiaoyan Lin, Zhi Tang, Xiaofan Lin, and Josef B. Baker. Wikimirs: A mathematical information retrieval system for wikipedia. *Proceedings of the 13th ACM/IEEE-CS joint conference on Digital libraries. Pages 11-20*, 2013.
- [14] David Stalnaker and Richard Zanibbi. Math expression retrieval using an inverted index over symbol pairs in math expressions: The tangent math search engine at ntcir 2014. *Proc. SPIE 9402, Document Recognition and Retrieval XXII*, 940207, 2015.
- [15] David Stalnaker and Richard Zanibbi. Math expression retrieval using an inverted index over symbol pairs. *Proc. SPIE 9402, Document Recognition and Retrieval XXII*, 940207, 2015.
- [16] Shahab Kamali and Frank Wm. Tompa. A new mathematics retrieval system. *CIKM*, 2010.
- [17] Kai Ma, Siu Cheung Hui, and Kuiyu Chang. Feature extraction and clustering-based retrieval for mathematical formulas. In *Software Engineering and Data Mining (SEDM), 2010 2nd International Conference on*, pages 372–377, June 2010.
- [18] Cyril Laitang, Mohand Boughanem, and Karen Pinel-Sauvagnat. Xml information retrieval through tree edit distance and structural summaries. In *Information Retrieval Technology*, volume 7097 of *Lecture Notes in Computer Science*, pages 73–83. Springer Berlin Heidelberg, 2011.
- [19] Shahab Kamali and FrankWm. Tompa. Structural similarity search for mathematics retrieval. In *Intelligent Computer Mathematics*, volume 7961 of *Lecture Notes in Computer Science*, pages 246–262. Springer Berlin Heidelberg, 2013.
- [20] Richard Zanibbi and Bo Yuan. Keyword and image-based retrieval for mathematical expressions. *Multi-disciplinary Trends in Artificial Intelligence. 6th International Workshop*, pages 23–34, 2011.
- [21] Li Yu and Richard Zanibbi. Math spotting: Retrieving math in technical documents using handwritten query images. *Document Analysis and Recognition (ICDAR)*, pages 446 – 451, 2009.
- [22] T. Nguyen, S. Hui, and K. Chang. A lattice-based approach for mathematical search using formal concept analysis. *Expert Systems with Applications*, 2012.
- [23] Hiroshi Ichikawa, Taiichi Hashimoto, Takenobu Tokunaga, and Hozumi Tanaka. New methods of retrieve sentences based on syntactic similarity. *IPSJ SIG Technical Reports*, pages 39–46, 2005.
- [24] Yoshinori Hijikata, Hideki Hashimoto, and Shogo Nishida. An investigation of index formats for the search of mathml objects. In *Web Intelligence/IAT Workshops*, pages 244–248. IEEE, 2007.
- [25] Yokoi Keisuke and Aizawa Akiko. An approach to similarity search for mathematical expressions using mathml. *Towards a Digital Mathematics Library. Grand Bend, Ontario, Canada*, pages 27–35, 2009.
- [26] Yoshinori Hijikata, Hideki Hashimoto, and Shogo Nishida. Search mathematical formulas by mathematical formulas. *Human Interface and the Management of Information. Designing Information, Symposium on Human Interface*, pages 404–411, 2009.
- [27] PhilipN. Klein. Computing the edit-distance between unrooted ordered trees. volume 1461 of *Lecture Notes in Computer Science*, pages 91–102. Springer Berlin Heidelberg, 1998.

ID	formula	ID	formula
1	$\int_0^\infty dx \int_x^\infty F(x,y)dy = \int_0^\infty dy \int_0^y F(x,y)dx$	2	$X(i\omega)$
3	$x^n + y^n = z^n$	4	$\int_{-\infty}^\infty e^{-x^2} dx$
5	$\frac{f(x+h)-f(x)}{h}$	6	$\frac{\sin x}{x}$
7	$ax^2 + bx + c$	8	$\frac{e^x+y}{z}$
9	$O(n \log n)$	10	$H^n(X) = Z^n(X)/B^n(X)$
11	$A_n = \frac{1}{\pi} \int_{-\pi}^\pi F(x) \cos(nx) dx$	12	$\lim_{x \rightarrow \infty} (1 + \frac{1}{x})^x$
13	$f(x) = f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \dots$	14	$f(a) = \frac{1}{2\pi i} \oint_r \frac{f(z)}{z-a} dz$
15	$x^2 + 2xy + y^2 = x ^2 + 2 x y + y ^2$	16	$\int_a^b f(x) dx = F(b) - F(a)$
17	$\frac{n!}{r_1! \cdot r_2! \cdot \dots \cdot r_k!}$	18	$-b \pm \sqrt{b^2 - 4ac}$
19	$1 + \tan^2 \theta = \sec^2 \theta$	20	$\vec{u} = (x, y, z)$

Table 1: Test query set

Query ID	Relevance Score					Total judged
	0	1	2	3	4	
1	15	2	2	0	1	20
2	19	0	0	0	0	19
3	15	4	1	0	0	20
4	0	0	0	6	14	20
5	0	0	0	8	12	20
6	0	2	5	2	10	19
7	1	4	3	3	9	20
8	5	1	1	1	0	8
9	0	0	4	11	5	20
10	0	0	1	16	3	20
11	0	0	0	1	6	7
12	0	0	4	13	3	20
13	14	3	1	1	1	20
14	0	2	5	8	5	20
15	0	0	0	15	5	20
16	8	6	2	2	2	20
17	0	0	5	13	2	20
18	19	1	0	0	0	20
19	19	1	0	0	0	20
20	19	1	0	0	0	20

asdfasdf

Query ID	Relevance Score					Total judged
	0	1	2	3	4	
1	13	2	2	1	1	19
2	15	1	0	3	1	20
3	0	0	0	0	20	20
4	0	0	0	1	19	20
5	0	0	0	0	20	20
6	1	0	0	0	19	20
7	0	0	0	0	20	20
8	4	0	1	2	0	7
9	0	2	5	8	5	20
10	0	0	0	9	11	20
11	0	0	0	2	5	7
12	0	0	12	3	5	20
13	15	1	2	1	1	20
14	0	0	0	0	20	20
15	0	0	0	13	6	19
16	0	0	2	0	18	20
17	0	0	0	0	20	20
18	0	0	2	7	11	20
19	0	0	2	7	11	20
20	0	0	2	7	11	20

asdfasdf

Table 2: Distribution table