

**A NOVEL SIMILARITY-SEARCH METHOD FOR MATHEMATICAL  
CONTENT IN  $\text{\LaTeX}$  MARKUP AND ITS IMPLEMENTATION**

by

Wei Zhong

A thesis submitted to the Faculty of the University of Delaware in partial fulfillment of the requirements for the degree of Master of Science in Electrical and Computer Engineering

Summer 2015

© 2015 Wei Zhong  
All Rights Reserved

**A NOVEL SIMILARITY-SEARCH METHOD FOR MATHEMATICAL  
CONTENT IN  $\text{\LaTeX}$  MARKUP AND ITS IMPLEMENTATION**

by

Wei Zhong

Approved: \_\_\_\_\_

Fouad E. Kiamilev, Ph.D.

Professor in charge of thesis on behalf of the Advisory Committee

Approved: \_\_\_\_\_

Kenneth E. Barner, Ph.D.

Chair of the Department of Electrical and Computer Engineering

Approved: \_\_\_\_\_

Babatunde A. Ogunnaike, Ph.D.

Dean of the College of Engineering

Approved: \_\_\_\_\_

James G. Richards, Ph.D.

Vice Provost for Graduate and Professional Education

## **ACKNOWLEDGMENTS**

Thank you to my family for their support from every perspective through out my graduate academic education. Thank you to my advisor Hui Fang who offers me the opportunity to develop my idea further and supports me in many other ways. I am also grateful to all InfoLab members for their kind help. And thanks to those not previously mentioned, who have influenced me or helped along the way.

# TABLE OF CONTENTS

|                                    |            |
|------------------------------------|------------|
| <b>ABSTRACT</b>                    | <b>vii</b> |
| <b>Chapter</b>                     |            |
| <b>1 BACKGROUND</b>                | <b>1</b>   |
| 1.1 Math IR Domains                | 1          |
| 1.2 Issues in Measuring Similarity | 3          |
| 1.3 Related Work                   | 5          |
| 1.3.1 Text-based methods           | 5          |
| 1.3.2 Structure-based methods      | 6          |
| 1.3.3 Other related work           | 8          |
| 1.3.4 Performance Review           | 9          |
| <b>2 METHODOLOGY</b>               | <b>10</b>  |
| 2.1 Intuitions                     | 10         |
| 2.1.1 Commutative immunity         | 11         |
| 2.1.2 Sub-structure query ability  | 11         |
| 2.1.3 Index and search properties  | 11         |
| 2.2 Structure Similarity           | 12         |
| 2.2.1 Definitions                  | 13         |
| 2.2.1.1 Formula tree               | 14         |
| 2.2.1.2 Formula subtree            | 14         |
| 2.2.1.3 Leaf-root path set         | 14         |
| 2.2.1.4 Index                      | 15         |
| 2.2.2 Search method                | 15         |

|          |                                                            |           |
|----------|------------------------------------------------------------|-----------|
| 2.2.3    | Substructure matching . . . . .                            | 16        |
| 2.2.3.1  | Observation #1 . . . . .                                   | 16        |
| 2.2.3.2  | Observation #2 . . . . .                                   | 17        |
| 2.2.3.3  | Observation #3 . . . . .                                   | 17        |
| 2.2.4    | Interpretation . . . . .                                   | 18        |
| 2.2.5    | The decompose-and-match algorithm . . . . .                | 20        |
| 2.3      | Symbolic Similarity . . . . .                              | 21        |
| 2.3.1    | Ranking constrains . . . . .                               | 23        |
| 2.3.2    | The mark-and-cross algorithm . . . . .                     | 24        |
| 2.4      | Combine the Two . . . . .                                  | 28        |
| 2.4.1    | Relaxed structure match . . . . .                          | 28        |
| 2.4.2    | Matching-depth . . . . .                                   | 29        |
| 2.4.3    | Matching-ratio . . . . .                                   | 30        |
| 2.4.4    | Final ranking schema . . . . .                             | 30        |
| 2.4.5    | Illustrated by an example . . . . .                        | 30        |
| <b>3</b> | <b>IMPLEMENTATION AND EVALUATION . . . . .</b>             | <b>35</b> |
| 3.1      | Implementation . . . . .                                   | 35        |
| 3.1.1    | Crawler . . . . .                                          | 35        |
| 3.1.2    | Parser . . . . .                                           | 35        |
|          | <b>REFERENCES . . . . .</b>                                | <b>40</b> |
|          | <b>Appendix</b>                                            |           |
|          | <b>A TITLE OF APPENDIX A . . . . .</b>                     | <b>44</b> |
|          | <b>B TITLE OF APPENDIX B . . . . .</b>                     | <b>45</b> |
|          | <b>List of Tables</b>                                      |           |
| 2.1      | First two iterations of example score evaluation . . . . . | 33        |
| 2.2      | 3rd iteration of example score evaluation . . . . .        | 33        |

|     |                                                                      |    |
|-----|----------------------------------------------------------------------|----|
| 2.3 | 4th iteration of example score evaluation . . . . .                  | 34 |
| 3.1 | Test query set . . . . .                                             | 36 |
| 3.2 | Query statistics for matching structurally only . . . . .            | 37 |
| 3.3 | Query statistics for matching structurally with symbolic scoring . . | 38 |

## List of Figures

|     |                                                            |    |
|-----|------------------------------------------------------------|----|
| 2.1 | Leaf-root path example . . . . .                           | 12 |
| 2.2 | Leaf-root paths with different structure . . . . .         | 19 |
| 2.3 | Formula subtree matching . . . . .                         | 20 |
| 2.4 | The decompose-and-match algorithm . . . . .                | 22 |
| 2.5 | The mark-and-cross algorithm . . . . .                     | 26 |
| 2.6 | Example query/document expression representation . . . . . | 31 |

## ABSTRACT

Mathematical content are widely contained by digital document, but major search engines fail to offer a way to search those structural content effectively, because traditional IR methods are deficient to capture some important aspects of math language. In this paper, we propose a similarity-search method for L<sup>A</sup>T<sub>E</sub>X math expressions trying to provide a new idea to better search math content. Our approach uses an intermediate tree representation to capture structural information of math expression, and based on a previous idea, we index math expressions by tree leaf-root paths. A search method to limit search set for possible subexpression isomorphism is provided. We rank search results by a few intuitive similarity metres from both structural and symbolic points of view. We also build our own proof-of-concept prototype search engine to demonstrate these ideas, and thus are able to present some evaluation results through this paper. Experiment shows these proposed measurements can advance effectiveness with respect to our baseline search method.

## Chapter 1

### BACKGROUND

Apart from general text content, structured information is also widely contained by digital document. Among these, a lot of mathematical content are represented, they are primarily using  $\text{\LaTeX}$ , which is a rich structural markup language. Information Retrieval on those structured data in mathematics language is not as well-studied or exhaustively covered as that is in general text IR research. To better search mathematical content can be significantly meaningful in terms of extending our border on information retrieval.

However, the structured sense of mathematical language, as well as many its semantic properties (see section 1.2), makes general text retrieval models deficient to provide good search results. This is because some fundamental differences between mathematical language and general text. Through this paper, we have made our efforts to tackle some of the problems we are having to search mathematical language. Some of the ideas used in this paper that deals with "tree structured" data can be generalized and potentially applied to other fields of structured data retrieval.

#### 1.1 Math IR Domains

Mathematical information involves a wide spectrum of topics, [1] gives a comprehensive review on mathematical IR researches. We are of course not focusing on every aspect in mathematical information retrieval. It is good to clarify our concentration in this paper here by first listing a set of concentrations in the related research area and define our target field of study.



Listed here, are considered three major topics for mathematical information retrieval:

1. Boolean or Similarity Search
2. Math Detection and Recognition
3. Evaluation, Derivation and Calculation

Boolean/Similarity search finds or ranks mathematical expressions against a query. They define the criteria of matching expressions or dimensions of similarity between two expressions. This is analogy to the boolean or similarity search of general text search engines, except the query and document may contain mathematical expressions. Some search engines deal with only formula (e.g. SearchOnMath <sup>1</sup>, Uniquation <sup>2</sup> and Tangent <sup>3</sup> ) and some math-aware search engines (e.g. WolframAlpha <sup>4</sup> and Zentralblatt math from MathWebSearch <sup>5</sup> ) are able to search query with mixed text and mathematical formula. These search engines can be useful in many ways, for example, student may utilize it to know which identity can be applied to a formulae in order to give a proof of that formulae. This is also the area where we focus in this paper.

Digital mathematical content document can also be in an image format (e.g. a handwritten formula), topics on retrieving information in these image requires detection and recognition of their visual features (texture, outline, shape etc.).

Also, because the nature of mathematical language, a query (e.g. an algebra expression) can potentially derived into alternate forms, or calculated and evaluated into a value. These potentially require a system to handle symbolic or numeric calculation, or even a good knowledge of derivation rules implied by different mathematical

---

<sup>1</sup> <http://searchonmath.com>

<sup>2</sup> <http://uniquation.com/en>

<sup>3</sup> <http://saskatoon.cs.rit.edu/tangent/random>

<sup>4</sup> <https://www.wolframalpha.com/>

<sup>5</sup> <http://search.mathweb.org/zbl/>

expression. Those numeric search engines (e.g. computational engine *Symbolab*<sup>6</sup> and WolframAlpha) can help evaluate mathematical expressions and reveal the deeper information contained from those expressions.

Besides the first three concentrations, there are many other topics. Knowledge mining, for example, will need deeper level of understanding on math language. A typical goal of this topic is to give a solution or answer based on all the document information retrieved. e.g. “Find an article related to the *Four Color Theorem*” [2].

These topics somehow overlap in some cases, for example, some derivation can be used to better assess the similarity between math formulae, e.g.  $\frac{a+b}{c}$  and  $\frac{a}{c} + \frac{b}{c}$  should be considered as relevant because their equivalent form of representation. Similarly, mathematical knowledge is required to understand the same meaning (thus high similarity) between  $\binom{n}{1}$  and  $C_n^1$ . So boolean or similarity search possibly involves a level of understanding on mathematical language. However, we are not going to include these problems into our research domain, instead, this paper addresses mathematical expression similarity from only structural and symbolic perspectives.

## 1.2 Issues in Measuring Similarity

Unlike general text content, mathematical language, by its nature, has many differences from other textual documents, there are a number of new problems in measuring mathematical expression similarity. Even without caring about the possible derivations and high level knowledge inference, there are still a set of new problems for measuring mathematical similarity.

Firstly, the differences between mathematical expressions should be captured in a cooperative manner in order to measure similarity, because only respecting symbolic information is not sufficient in mathematical language. To illustrate this point, we know that  $ax + (b+c)$  in most cases is not equivalent to  $(a+b)x + c$  although they have the same set of symbols, because the different structure represents a different semantic meaning. And the order of tokens in math expression can be commutative in some

---

<sup>6</sup> Symbolab Web Search: <http://www.symbolab.com>

cases but not always, for example, commutative property in math makes  $a+b = b+c$  for addition operation, but on the other hand  $\frac{a}{b}$  is most likely not equivalent to  $\frac{b}{a}$ . These characteristics make many general text search methods (e.g. *bag of words* model, *tf-idf* weighting) inadequate. Moreover, symbols can be used interchangeably to represent the same meaning, e.g.  $a^2 + b^2 = c^2$  and  $x^2 + y^2 = z^2$ . However, interchanging symbols does not always preserve mathematical semantics, changes of symbols in expression preserve more syntactic similarity when changes are made by substitution, e.g. Given query  $x(1+x)$ , expression  $a(1+a)$  are considered more relevant than  $a(1+b)$ .

Secondly, how we evaluate structural similarity between expressions is a question. A complete query may structurally be a part of a document, or only some parts of a query match somewhere in a document expression. In cases when a set of matches occur within some measure of “distance”, we may consider them to contribute similarity as a whole, but when matches occur “far away” for a query expression, then under the semantic implication of mathematics, they probably will not contribute the similarity degree in any way. We need metrics to score these similarity under certain rules for relevance assessments.

Lastly, when trying to capture more semantic information from expressions, we can improve our measurement on similarity but it may introduce more ambiguity. For example, semantically incorrect math markups in document, e.g. using “sin” in L<sup>A</sup>T<sub>E</sub>X markup instead of macro “\sin”, will make it difficult to tell whether it is a product of three symbols or a *sine* function if we want to capture their semantical meaning in such a depth. And depending on what level of semantical meaning we want to capture, ambiguity cases can be different. Consider  $f(2x+1)$ , if we want to know if  $f$  is a function rather than a variable, the only possibility is looking for its contexts, but we can nevertheless always think of it as a product without losing the possibility to search similar expression like  $f(1+2y)$ , the same way goes reciprocal  $a^{-1}$  and inverse function  $f^{-1}$ . Most often, even if no semantic ambiguity occurs, efforts are needed to capture some semantical meanings. e.g. In  $\sin 2\pi$ , it is not easy to figure out, without a knowledge on trigonometric function convention, that  $2\pi$  is the subordivative of sine

function.

### 1.3 Related Work

Boolean or similarity 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]. A variety of approaches have already emerged in an early timeline [4], but we can nevertheless categorize them into several different ideas. [5, 6, 7] use the same way to classify them into text-based and tree-based (structure-based), here we decide to follow the same classification method and give an overview on those different ideas.

#### 1.3.1 Text-based methods

Many researchers are utilizing existing models to deal with mathematical search, and use text-based approaches to capture structural information on top of matured text search engine and tools (such as *Apache Lucene*).

DLMF project from NIST [8] uses “flattening process” to convert math to textualized terms, and normalize them into *sorted parse tree normal form* which creates an unique form for all possible orders of nodes among associative and commutative operators. Then further introduces serialization and scoping to stack terms [9], trying to capture structure information by using text-IR based systems that supports phrase search. Similar idea is also used by [10], expressions are also augmented for various possible representations, but variables are also replaced and normalized, but they are using postfix notation, allows to search subexpressions without knowing the operator between them. MIaS system [11, 12, 13], like the methods above, are also trying to reorder commutative operations, normalize variables and constant numbers into unified symbols, doing augmentation in a similar fashion. It indexes expressions and subexpressions from all depth levels. The system is able to discriminate and assign different weight based on their generalization level, and place emphasis on that a match in a complex expression is assigned higher weight [13].

Augmentation usually consists a storage demand for combination of both symbols (e.g.  $a$  and  $b$ ) and unified items ( $id$ ,  $const$ ) in different levels, in order to capture both symbolic information and structure information. Thus implies complex expressions with many commutative operators will cost inevitably larger storage space, the benefits of capturing expression variances will be overshadowed.

Although named as structured-based approach, [14] is using *longest common subsequence* algorithm to capture structure information (in an unified *preprocessed string* and a *level string*). The method takes  $O(n^2)$  complexity for comparing a pair of formulae, and no index method is proposed. Therefore is not feasible to efficiently apply to a large collection. The Mathdex search engine [15], from another perspective, uses query likelihood approach [16] to estimate how likely the document will generate the query expression. Math GO! [17] is another system advances some transitional method to better search math content. It tries to find all the symbols and map formula pattern to pattern name keywords (like *matrix* or *root*), and proposes to replace term frequency by co-occurrence of a term with other terms.

### 1.3.2 Structure-based methods

What text-based methods share in common is they are converting math language symbols to linear tokens, the intrinsic defect when using a bag of words to replace structured information will make conversion process lose considerable information or cause storage-inefficiency. In order to cope with the problems from text-based approach, structure-based methods generally generate intermediate tree structure, and use these information to index or search.

#### Term Indexing

Whelp [18] and MathWebSearch (MWS) directed by Kohlhase [19, 20, 21], are one of the notable structure-based ones which are derived from *automatic theorem proving* and *unification theory* [22]. The system of MWS uses *term indexing* [23] in a

*substitution tree index* [23] to minimize access time and storage. Because the subexpression is not easy to search using substitution tree, MWS indexes all sub-terms, but the increased index size remains manageable [19]. However, their index relies heavily on RAM memory, the considerable RAM resource usage (170GB reportedly [21]) makes it have to scale to accommodate 72% papers on arXiv.

### Leaf-root path

[24] uses leaf to root XML path in a MathML object to represent math formula. When efficiency is considered, it only indexes the first and deepest path (to indicate how a formula is started and presumably the most characteristic part of a formula); when user wants to obtain the perfect-match result, it indexes all the MathML object leaf-root path. The boolean search is performed by giving all the paths match with those of the search query. [25] further develops with incorporation of previous method using breath-first search, to add sibling nodes information into index and have achieved better effectiveness.

Very similar idea is proposed by [26] and used in [27]. The latter transform MathML to an “apply free” markup from which the leaf-root path are extracted. Leaf-root path is also used to evaluate similarity between two expressions in MathML.

### Symbol layout tree

A *symbol layout tree* [28] (SLT) or *presentation tree* [5] describes geometric layouts of mathematical expression. WikiMirs [5] uses two templates to parse L<sup>A</sup>T<sub>E</sub>X markup with two typical operator terms: explicit ones (“\frac”, “\sqrt”, etc.) and implicit ones (“+”, “÷”, 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. [28] uses symbol layout tree as a kind of substitution tree, each branch in the tree represents a binding chain for variables, and every child node is an instance of its parent for a generalized term. They introduce *baseline size* to help group similar expressions together in their substitution SLT in order to decrease

tree branch factor, however, this makes a single substitution variable not able to match multiple terms along the baseline. Also their implementation makes it unable to index certain formula (e.g. a left-side superscript) and have to generate many queries (e.g. all possible format variations and sub-expressions etc.) for a single query at the time of search. Later [6, 29] have developed a *symbol pairs* idea to capture a relative position information between symbol pairs. Due to the many possible combinations of symbol pairs in a complex math expression, the storage cost is intrinsically large. However, the key-value storage style will be suitable for fast lookup (e.g. HASH).

### Other structure-based methods

A novel indexing scheme and lookup algorithm is proposed by [30], its index has hash signature for each subtree because they have observed a lot of common subtree structure occur in math formula collection. This idea will result in a slower index growth. Their lookup algorithm supports wildcards, and performs a boolean match test. Although their lookup time is generally below 700ms, the index size where query lookup time is tested is unclear, but presumably no greater than 70,000 expressions. By constructing a DOM tree, [31] extracts semantic keywords, structure descriptions to indicate subordinative relationship in a string format. The similarity is calculated using normalized tf-idf vector (trained by clustering algorithms) by dot product. Although the final index is generated from text, promising results have been achieved. Tree edit distance is adopted by [32], it tries to overcome the bad time complexity of original algorithm by summarizing and using a structure-preserving compromised edit distance algorithm. Although the result shows query processing time is long but it is reduced to average 0.8 seconds by applying with an early termination algorithm along with a distance cache [33].

#### 1.3.3 Other related work

There are a number of articles trying to use image to assess math expression similarity. [34] compares their image-based approach using connected component-based

feature vector with a proposed text-based method, reported precision@k values are low, but the potential for this method to be combined with shape representations or other features will potentially improve it and make it valuable for searching mathematical expressions in image format. [35] uses X-Y tree to cut the page in vertical and horizontal directions alternatively, in order to retrieve math symbols from images.

A lattice-based approach [36] builds formal concept based on selected feature sets of each formula. The ranking is calculated by the distances from query in the lattice map when the query is inserted.

### 1.3.4 Performance Review

In the review of many related past research in section 1.3, we find the combination of state-of-art general text search engine (i.e. *Apache Lucene*) with the efforts to augment expressions by permutation and unification to satisfy the needs of mathematical search have achieved a good result in different metrics of evaluation: The text-based system of MlaS over-performs those of structure-based and ranked the first in four out of six measurement in NTCIR-11 conference [37]. However, structure-based method such as symbol pairs proposed in [6, 29] also gets a competitive result [37].

Although text-based methods have achieved relatively better effectiveness, their commonly adoption of augmentation makes it expensive in terms of hardware storage. On the other hand, structure-based method has the merit to best capture semantic information of a typical mathematical expression, thus still has the potential to further improve effectiveness given the fact that the tools and theory behind text-based methods are already mature and have been developed for decades.



## Chapter 2

### METHODOLOGY

Our method can be seen as an approach built upon the idea of leaf-root path or sub-path [26, 24, 27, 25] from an operation tree [1]. But we have developed this idea further in many ways. Our index is composed from leaf-root paths extracted from an operation tree, we simultaneously search along the way of all leaf-root paths from a given query operation tree, which is essentially traversing “reversed” sub-paths. Searching in this way also makes a pruning method possible to limit our search set. We also propose a sub-structure testing algorithm, which are utilizing some observed properties from our indexed tree. Apart from these, we provide a symbolic similarity algorithm to rank  $\alpha$ -equivalent expressions higher. The methods in a nutshell is, for a document expression, construct operation tree, break it down into sub-paths, and index those paths. For a query, traversal index tree as the same way of going through the reversed sub-paths of that query operation tree, search the sub-paths and their merged paths in index meanwhile calculate indexed expression similarity to rank each document expression.

#### 2.1 Intuitions

First it is beneficial to document our intuitions on using operation tree as our intermediate representation and our idea to index it in a way of reversed sub-path tree, and also explain in abstract why this way helps reduce index space and boost search speed. We will give an illustrative example to describe these processes further in section 2.4.5.

### 2.1.1 Commutative immunity

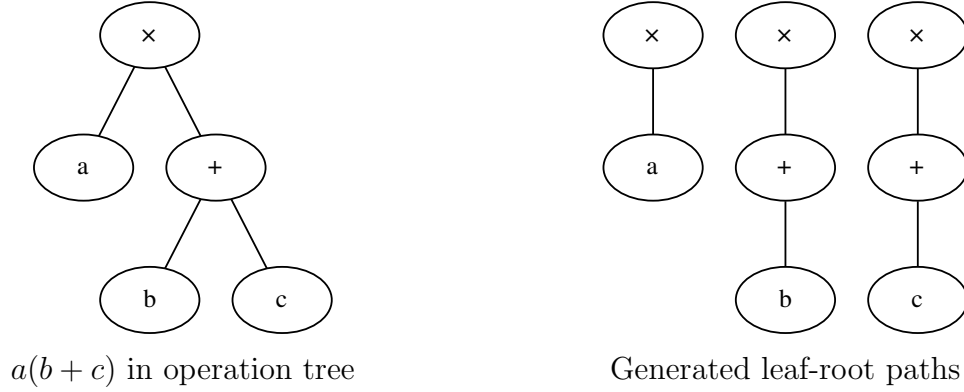
Operators with semantic implication of commutative property (e.g. addition and multiplication) are exhaustively used in mathematical language. The ability to identify the identical equations for any permutation is very essential for a mathematical similarity search engine. Given this as a starting point, the leaf-root paths have the advantage to cope this so that we do not need to generate different order of patterns to match formulae with commutative operator. To illustrate this, we know that a leaf-root path from an operation tree (see figure 2.1) is generated through traversing in a bottom-up (or top-down) fashion from a tree, therefore these sub-paths are independent to the relative position of commutative operands. In another word, an operation tree uniquely determines the leaf-node paths decomposed from the tree, no matter how operands are ordered.

### 2.1.2 Sub-structure query ability

On the other hand, the structure of operation tree also makes it easy to represent sub-expression relation with a formula, because a sub-expression in a formula is usually (depending on the way we construct an operation tree) also a subtree in an operation tree. And by going up from leaves of operation tree, we are essentially traversing to an expression from its subexpression for every level. By making all the leaf-root paths as an index, we can search an expression by going through and beyond the leaf-root paths from its subexpression. This makes operation tree better in terms of searching an expression given a sub-expression of it in query. And it avoids information augmentation on index as some other structure-based methods need (e.g. index all sub-terms of an expression in MWS [19]). Therefore it also helps save storage space.

### 2.1.3 Index and search properties

Additionally, some properties from this method suggest some reduce of space and pruning possibilities in search process. First the sub-paths themselves can be indexed into a tree so that we can search a sub-path by traversing a sub-path tree,



**Figure 2.1:** Leaf-root path example

instead of hashing it to find a corresponding value as the symbol-pair search engine (i.e. *Tangent* [6]) does. This allows us to save a lot of space as the reverted sub-paths of a large collection will have a great percentage of level sharing a common substring with each other. Also the way to search in a tree structure with a limited branch factor does not lose much efficiency compared to the HASH methods used in *Tangent*, while also offering great storage efficiency. Second, by searching from all the “reverted” sub-paths of a query expression in our proposed index, and applying an intersection on the results from different sub-paths, we will find all the expressions having that query as subexpression (number four observed property from section 2.2.3). And during this search process, upon going further from the query expression root in the “reverted” sub-path, we can merge the next search directories by pruning all the entries that are not shared in common among all the “reverted” search paths. Furthermore, if a relevant indexed formula has been found in a search path level, then other relevant sub-paths will be found at the same level too, thus some implementation strategies can be applied to speed search further (i.e. distributed search to quick search massive).

## 2.2 Structure Similarity

The basic ideas used in our approach, to test whether a mathematical expression is a sub-structure of another, to prune and constrain search set are the foundation

work in our research. It is desired to give a description in a formal language so that we can deliver these ideas in the most precise way. Some important observations as well as brief justifications are provided for completeness.

### 2.2.1 Definitions

For the second issue addressed in section 1.2, specifically, to assess the structural similarity. Some formal definitions are proposed from previous study [38] which provides a quantified  $n$ -similarity relation to address the similarity degree between math expressions, based on the max-weight common subtree. The subtree, by their definition, includes all descendants from a node, however, we are going to use the subtree definition in graph theory here to describe the sub-structure relation. To be explicit, given a rooted tree  $T$ , the connected graph whose edges are also in  $T$  is defined as the subtree of  $T$ . We are also going to give a definition to describe a “match” between query and a document math formula.

Because we are using some formal descriptions to better illustrate our ideas, there are some notations we need to clarify that are used throughout this paper.

A *path*  $p$  is a sequence of numbers given by  $p = p_0p_1 \dots p_n$  where  $n \geq 0$ ,  $p_i \in \mathbf{R}$ .  $\mathbf{P}$  is the set of all paths. A *leaf-root path* is a path from root to a leaf in a tree. Any function  $f : \mathbf{R} \rightarrow \mathbf{R}$  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  ${}^1p = p_0p_1 \dots p_n$  and  ${}^2p = p_np_{n+1} \dots p_m$  where  $m \geq n$ , to be a new path denoted as  ${}^1p \cdot {}^2p = p_0p_1 \dots p_np_{n+1} \dots p_m$ , and the concatenation of a path  $p$  on a path set  $S = \{s_1, s_2 \dots s_n\}$  is defined as  $S \cdot p = \{s_1 \cdot p, s_2 \cdot p \dots s_n \cdot p\}$ . The *longest common prefix* path  $p^*$  between two paths  $p_1$  and  $p_2$  is mapped by the function named lcp, which is defined by  $p^* = \text{lcp}(p_1, p_2) = \text{lcp}(p_2, p_1)$ .

Each node in a formula tree is associated with a label (labels are not required to be distinct here) to represent an unified mathematical operator, variable, constance etc., a class of similar operators can have the same label value (e.g. same label for token  $+$ ,  $\oplus$  and  $\pm$ ). Also each leaf node is associated with a symbol value to represent a symbolic instance of that constance or variable (e.g. “123”,  $\beta$ ,  $x$ ,  $y$  etc.). Besides, a

*formula subtree* relation is also defined to address the sub-structure relation between two mathematical expressions. Below are our formal definitions.

#### 2.2.1.1 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  $\ell_T$ , and each leaf  $l \in V(T)$  is further associated with a symbol  $\mathcal{S}_T(l) \in \mathbf{R}$  mapped by symbol function  $\mathcal{S}_T$ . For convenience, we will write function  $\ell$  and  $\mathcal{S}$  as their short names which refer to the tree implied by the context, and we use  $\mathcal{S}(p)$  to indicate the symbol of the leaf in a leaf-root path  $p$ . In addition, we use  ${}^sT$  to denote a subtree in  $T$  rooted by vertices  $s \in V(T)$ , with all its descendants.

#### 2.2.1.2 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  $T$ .

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

#### 2.2.1.3 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  $g(T)$ .

#### 2.2.1.4 Index

An *index*  $\Pi$  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  $\Pi$  and  $\mathcal{I}_\Pi$  is called *index look-up function* for index  $\Pi$ .

#### 2.2.2 Search method

For a collection of document expressions, we will index them by merging all the reverted leaf-root paths from each document formula tree into a large “inverted” index tree, in which each node at path  $a$  stores the information of all the indexed formula trees in  $\mathcal{I}_\Pi(a)$ . Through searching all sub-paths from a query formula tree  $T_q$  at the same time, we are able to limit the set of possible formula trees being structurally matching (in formula subtree relation) with  $T_q$ , to only a subset of our index. This is illustrated as follows.

Given an index  $\Pi$  and a formula tree  $T_q$ ,  $\forall T_d \in \Pi$ : If  $T_q \preceq_l T_d$  on  $\Phi$ , 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 2.2.1.2 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  $\Pi$  with respect to each of the elements in  $L$ , that is to say  $\forall a \in L$ ,  $T_d \in \mathcal{I}_\Pi(a)$ .

To summarize, 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 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  $\Pi$ .

### 2.2.3 Substructure matching

However, query formula tree will not necessarily be formula subtree of 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 2.2. To address this problem, we propose an algorithm described in figure 2.4, to test the document formula trees in our search set to see if they are in formula subtree relation with query formula tree. This algorithm is inspired from the following observations.

#### 2.2.3.1 Observation #1

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

$$\deg(v) \leq \deg(\phi(v)) \quad (2.1)$$

$$\ell(p) = \ell(\phi(p)) \quad (2.2)$$

$$|g(T_q)| \leq |g(T_d)| \quad (2.3)$$

*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 2.2.1.2. Therefore any different edge from  $v$  is associated with a distinct edge from  $\phi(v)$ , thus we can get (2.1). 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 (2.2) 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 (2.3) can be obtained in a similar manner with that of (2.1), except neighbor edges are replaced by leaf-node paths.

### 2.2.3.2 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 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. Since  $\ell(p) \subseteq \ell(g(T_d))$ , we know that there must exist a path  $p' = p'_0 \dots p'_n \in g(T_d)$  such that  $\ell(p) = \ell(p')$  where  $p'_n$  is the leaf of  $T_d$ . Then 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  $T_d$ .

### 2.2.3.3 Observation #3

For two formula trees  $T_q$  and  $T_d$ , if  $T_q = T(V, E, r) \preceq_l T_d$  on  $\Phi$ ,  $\forall a, b \in g(T_q)$  and a mapping  $\phi \in \Phi$ . Let  $T'_d = {}^tT_d$  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  $b = b_0 \dots b_n b_{n+1} \dots b_{m-1}$ , 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  $b' \in g({}^tT_d)$ , therefore  $a'_0 = \phi(a_0) = \phi(r) = t = b'_0$ . For the first conclusion, if  $b' = \phi(b)$ , there are two cases. If either  $|a|$  or  $|b|$  is equal to one then  $|\text{lcp}(a, b)| = |\text{lcp}(a', b')| = 1$ ; 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 both 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,



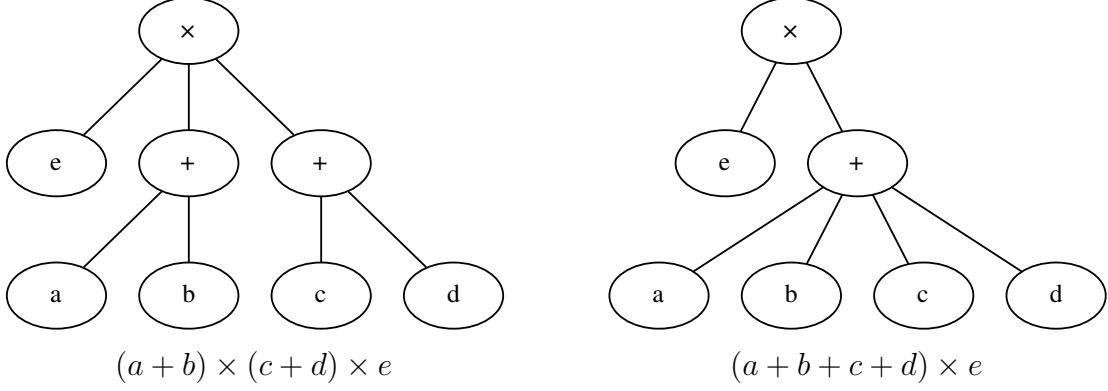
because  $|\text{lcp}(a, c)| \neq |\text{lcp}(a, b)| = |\text{lcp}(a', b')|$ , thus  $|\text{lcp}(a, c)| \neq |\text{lcp}(a', b')|$  which is impossible.

#### 2.2.4 Interpretation

The observations above offer some insights on how to test a substructure of a mathematical expression.

First we give some explanations on the definition. A formula subtree relation defined in 2.2.1.2 describes not only a sub-structure relation between two math expressions, it also requires a label similarity and leaf inclusion. Because structure shape (subtree isomorphic) is not only one factor to determine whether a math formula is a subexpression of another. Given expression in figure 2.1 as an example,  $b + c$  and  $a \times (b + c)$  are considered “similar” because  $b + c$  is structurally an subexpression of  $a \times (b + c)$ . However, if expressions with different symbols but in similar semantics are given, e.g.  $b \oplus c$  or  $b \pm c$ , they should also be considered as similar to  $a \times (b + c)$  because both the operations has the similar semantical meaning as “add”. These operations should be labeled in which all the similar operation tokens have the same label value. Also, operation tree representation generally puts operator in the intermediate nodes and operands in the leaves, so it is not common to address a sub-structure without leaves, like “ $a \times +$ ”. This is the reason that a structure-similarity relation of two should also contain their leaves.

Now that we have defined our structure similarity rule as whether two trees  $T_q$  and  $T_d$  can satisfy:  $T_q \preceq_l T_d$ , we break down a formula tree into leaf-root paths  $p$  and index the label of each path  $\ell(p)$ . So if given a “similar” path  $q$ , we can further find the previous trees that also have  $\ell(q)$  as its labeled path. In section 2.2.3, the first observation gives some constrains to test if two leaf-node paths are similar without knowing the complete tree from which they are generated. However, comparing all the paths from the index one by one would be very inefficient. Section 2.2.2 suggests if we search the index by all the generated leaf-node paths from a tree at the same time, then we may just need to look into an intersected region instead of the whole collection.

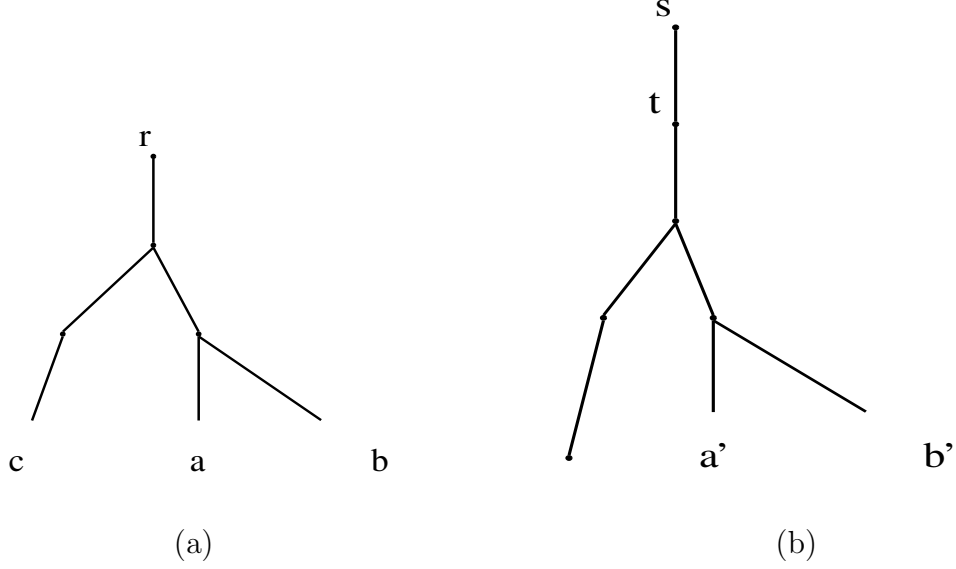


**Figure 2.2:** Leaf-root paths with different structure

Because every tree indexed ( $T_d \in \Pi$ ) and matched by the query will be found at our defined search set.

However, knowing the matched tree in in a set does not necessarily mean all the tree in the set match with query. Figure 2.2 gives one case where two set of leaf-root paths are identical while the structures from which they are generated are different, and not in any sub-tree relation. If left figure is the query, then we can certainly find the tree in right figure as long as it is indexed, indicated by section 2.2.2. Although leaf-root paths offer some desired properties, whether the trees found through searching sub-paths of a query are also structure isomorphic with the query tree is still unknown.

Observations #2 and #3 in section 2.2.3 offer the way to test structure isomorphic. The former is a sufficient condition to test structure isomorphic, but the tree must first have only one leaf-root path. The latter states two necessary conditions to be a formula subtree of another. This leads to an idea to decompose the tree and divide the problem into subtree matching problems by ruling out impossible matches between leaf-root paths using observation #3, until it is obvious to conclude the structure isomorphic in a sub-problem by using observation #2. These observations inspire the idea to filter expressions with structure isomorphism in our search set  $\bigcap_{a \in L} \mathcal{I}_{\Pi}(a)$ .



**Figure 2.3:** Formula subtree matching

### 2.2.5 The decompose-and-match algorithm

Here we propose and describe an algorithm for formula subtree matching based on the interpretation in section 2.2.4. Figure 2.3 illustrates a general case in which query tree (a) is trying to match a document tree (b).

Initially every leaf-root path in (a) should be associated with a set of leaf-root paths in (b) that are possible (by the constraints of observation #1 in 2.2.3) to be isomorphic, we call this set *candidate set*. Given figure 2.3, the candidate set of path  $a$  in (a) probably is  $\{a', b'\}$  in (b) if the nodes are assigned universally the same label. Then we arbitrarily choose a path in (a) as a reference path (heuristically a *heavy path* [39]), 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 constraints from observation #3 and ruling out some impossible isomorphic paths in candidate set of each path in (a), then 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  $c$  anymore 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 very simple subproblem is encountered (e.g. the case in observation #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 detailed algorithm is described in figure 2.4. The main procedure is *decomposeAndMatch* where the parameters  $Q$  and  $C$  are 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 match is found, otherwise FAIL is returned indicating the formula tree in (a) cannot be a formula subtree of that in (b).

### 2.3 Symbolic Similarity

Besides structural similarity, symbolic similarity is also essential to be considered, because it has some benefits to further differentiate math expression similarity.

Firstly, it is essential to score equations with symbolic matches a higher similarity because they may imply more semantic meaning. For example,  $E = mc^2$  is considered more meaningful when exact symbols are used rather than just being structure identical with  $y = ax^2$ .

Secondly, as illustrated in section 1.2, same mathematical symbols in an expression (or bound variables) usually can only maintain semantical equality if the changes are made by substitutions. (similar to the notion of  $\alpha$ -equality [40]). This is an important semantic information that we need to capture in most math expressions that contain bound variables, to capture this kind of equality between expressions certainly involves comparison of symbols.

There is one thing to address here, in many mathematical search systems, a query may be specified with wildcards and thus will match any document with an expression substitution to that wildcard. And a query symbol not specified by wildcard

```

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 DECOMPOSEANDMATCH( $Q_t, C$ ) = FAIL then
22:          return FAIL
23:  return SUCC
24:
25: procedure DECOMPOSEANDMATCH( $Q, C$ )
26:  if  $Q = \emptyset$  then return SUCC
27:   $a := \text{OnePathIn}(Q)$  ▷ 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 MATCH( $a, a', Q_{\text{new}}, C_{\text{new}}$ ) then return SUCC
33:  return FAIL

```

**Figure 2.4:** The decompose-and-match algorithm

is expecting an exact symbolic occurrence in document. Yet in our approach, same symbols imply a higher rank, we can get an exact symbolic match if it does exist in our index. On the other hand, we do not have to specify wildcard to match its  $\alpha$ -equalities. [11] also doubts the wildcard demand in mathematical search as it is not common to expect an exact symbol occurrence when we query in mathematical language.

### 2.3.1 Ranking constrains

As we have discussed, symbolic similarity is essential to be captured, in order to further rank document expressions. Here we propose two constrains to addressed all the considerations, they can be summarized as:

- Given two formula trees  $T_q \preceq_l T_d$  on  $\Phi$ , suppose a leaf  $l \in V(T_q)$  is isomorphic to leaf  $l' \in V(T_d)$ , that is to say, there exists  $\phi(l) = l'$ , where  $\phi \in \Phi$ , then if their symbol matches, i.e.  $\mathcal{S}(l) = \mathcal{S}(l')$ , we score them higher than those do not match symbolically. And the more symbolic matches there are, the higher symbolic relevance degree two expressions would expect to have.
- $\alpha$ -equivalent expressions have more symbolic relevance degree than those are not, and the more *bound variables* (variable with same symbol in an expression) two expressions match at the structurally matching positions, the more symbolically relevant they are considered to be.

In this paper, we use these two constrains as the rule to rank retrieval results in terms of symbolic similarity. And in cases where both constrains can be applied, we prioritize the second constrain. This is because, intuitively, as long as the semantic meaning of two expressions is the same, using different set of symbols does not make a difference. However, bound variable match is more important because an mathematical expression with more than one identical symbols most often implies those symbols represents the same variable.

The constrains and idea above are illustrated by the following example. Let the rank of a document math expression  $d$  be  $r(d)$ , and given query  $\sqrt{a}(a-b)$  for instance. It is easy to see under the first constrain:

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

The document with the highest rank here is an exact match, with three symbols matching in total. The second document has two symbols matching while the third document has no symbolic match at all.

In the same manner, by the second constrain we have:

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

The first document uses bound variable  $x$  but it preserves the same semantics compared to the query, except for the symbol of that bound variable is different with that of query. The second one does not have bound variable match, in another word, it uses different symbols (i.e. “ $x$ ,  $y$ ”) at positions where query expression uses the same symbols (i.e. “ $a$ ”).

One common pattern the first example follows is all “ranking” documents in the first example have the bound variable match, the only difference is the number of symbol matches. On the other hand, all “ranking” documents in the second example have the same number of symbol matches (only “ $b$ ” is matched in a symbolic way). So it is easy to follow only one of the two constrain. However, sometimes both constrains can be applied so that conflict may occur and we have to choose only one to follow. For example, given document expression  $\sqrt{a}(x - b)$  and  $\sqrt{x}(x - b)$ , the former has two symbolic matches (i.e. “ $a$ ,  $b$ ”) with our query  $\sqrt{a}(a - b)$ , while it does not have bound variable match. The latter, on the other hand, has bound variable match while only has one symbolic match (i.e. “ $b$ ”). We nevertheless score the latter higher because it does not lose any semantics.

### 2.3.2 The mark-and-cross algorithm

Here we propose an algorithm to score symbolic similarity between query and document expressions. To follow all the constrains and issues addressed, intuitively, we first take the bound variable with greatest number of occurrence in query expression, try to match as much as possible with each bound variable from document expression. The *best-matching* bound variable in document expression is chosen to contribute to the

final symbolic relevance score (proportionally to the number of matches in that bond variable), and we exclude its paths from matching query paths in future iterations. In the next iteration, we choose the bond variable with the second number of occurrence in query expression and repeat this process until all the query bond variables are iterated.

The algorithm is described in figure 2.5. It 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 defined by a set  $V(D) = \{x \mid \mathcal{S}(x), x \in D\}$ , which contains all the leaf node symbols from document expression. Procedure *sortBySymbolAndOccur* takes a set of leaf-root paths and return a list of all the leaf-root paths. The list is sorted by tuple  $(N_p, \mathcal{S}(p))$  for list element  $p$ , where  $N_p$  is the number of  $\mathcal{S}(p)$  occurred in all the list path symbols. Take the example expression  $\left(b \cdot \frac{a+b}{a+c} + a\right)$  again, the result list returned by *sortBySymbolAndOccur* is a sequence of leaf-root path of  $a, a, a, b, b, c$  respectively.

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 function  $\text{sim}(a, a')$  measures the symbolic 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 similarity weights to leaf-root paths with exact symbol match. (this function  $\text{sim}(a, a')$  is used only when a document path is in candidate set, see figure 2.5).

Let us determine the proper value for  $\alpha$ . Consider the conflicting cases stated in this section by using another example here, given query expression  $a + \frac{1}{a} + \sqrt{a}$  and



```

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 := SORTBYOCCURANDSYMBOL( $Q$ )
10:  for  $a$  in QList do
11:    for  $v \in \mathcal{V}(D)$  do
12:       $m := -\infty$ 
13:       $m_p := \emptyset$ 
14:      for  $a' \in C_a \cap \{y \mid \mathcal{S}(y) = v, y \in D\}$  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  $a' \in C_a \cap \{y \mid \mathcal{S}(y) = v, y \in D\}$  do
38:          if  $T_{a'} = \text{mark}$  then
39:             $T_{a'} := \text{nextState}$ 
40:  return score

```

**Figure 2.5:** The mark-and-cross algorithm

document expression  $a + \frac{1}{a} + b + \frac{1}{b} + \sqrt{b}$ , we consider bond-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 latter. That is to say, given a document bond-variable matching  $k$  variables with that in query, we need  $\alpha$  to satisfy  $k\alpha > (k-1) \times 1 = k-1$  and  $\alpha < 1$ . Therefore, in our practice, we set  $\alpha$  to a value close to 1, e.g. 0.9.

By sorting the query paths in  $Q$ , the algorithm is able to take out paths from same bond 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 this value on  $B_v$  indicating the similarity between currently evaluating query bond variable and the bond variable  $v$  in document expression. 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), let the document bond variable  $m_v$  with greatest  $B_v$  value  $m$ , be the best-matching bond variable in  $\mathcal{V}(D)$  with the query bond variable just iterated, then add  $m$  to the 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 paths

are iterated, finally return the score indicating the symbolic similarity between the two expressions.

## 2.4 Combine the Two

We have already discussed the problem and proposed the algorithms to both test structure isomorphic and measure symbolic similarity between mathematical expressions. The two algorithms, however, do not cooperate in an unified way. To illustrate this point, assume we first use DECOMPOSE-AND-MATCH algorithm and have concluded one is a formula subtree of another, but we only get one possible substructure match. There are very likely other possible positions where this tree can also be formula subtree of another, because the candidate sets are not unique (or  $|\Phi| > 1$ ). Thus we need to exhaust all possibilities and apply MARK-AND-CROSS algorithm to each of them, in order to find the maximum symbolic similarity pair. Analogously, assume we first want to get the symbolic similarity degree, then we are uncertain about the paths we have specified in candidate set are really isomorphic paths to the query path, i.e. we will not guarantee substructure relation when using DECOMPOSE-AND-MATCH algorithm.

One possible way to both test strict structure isomorphism and measure symbolic similarity is to decompose the tree and try to match all the substructures but at the same time heuristically choose reference paths to find the best symbolic match in a greedy way, if no possible substructure can be matched isomorphically, we have to rollback and try other candidates using backtracking. From efficiency perspective, the MARK-AND-CROSS algorithm already has an worst case time complexity of  $O(|Q| \times |D|)$ , if we try to find the best match (find the one with not only the most symbolic similarity but also who satisfies structure isomorphism), then it will eventually lead to more time complexity.

### 2.4.1 Relaxed structure match

The complexity introduced to combine the two methods will make our approach infeasible to efficiently deal with large data set. Here we choose to relax our constrain

on strict structure isomorphism. As figure 2.2 has illustrated, we know that the labels of a leaf-node path set being a subset of another does not necessarily mean the tree generating the former path set is a formula subtree of that generating the latter. To further generalize it, we say 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 considering those cases will offset the benefit to strictly identify the structure isomorphism. Therefore, we loose our constrain on structure similarity so that any  $T_d \in \bigcap_{a \in L} \mathcal{I}_\Pi(a)$  is considered structurally relevant to query formula tree  $T_q$  in section 2.2.2, and we say  $T_d$  is *searchable* by  $T_q$  in index  $\Pi$ . Optionally, we can apply constrain 2.1 in observation #1 to further eliminate cases such as the one in figure 2.2 so that less query/document expressions which are not in formula subtree relation would be considered relevant. We name this constrain as *fan-number constrain*.

The revised method will collect all the “structurally similar” document formula trees that satisfy the fan-number constrain in the set  $\bigcap_{a \in L} \mathcal{I}_\Pi(a)$ . Then use MARK-AND-CROSS algorithm to rank their symbolic similarity.

### 2.4.2 Matching-depth

On the other hand, we introduce a *matching depth*  $d = |\hat{a}|$  into symbolic similarity measurement algorithm, to make it also consider the depth of a structural match. As it is addressed in [11], the deeper level at where two formula matches, the lower similarity weight would it be, since 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 impact on score, we define *matching depth factor*  $f(d)$  to be a function value where  $f$  is in negative correlation with matching depth:  $f(d) = 1/(1 + d)$ .

### 2.4.3 Matching-ratio

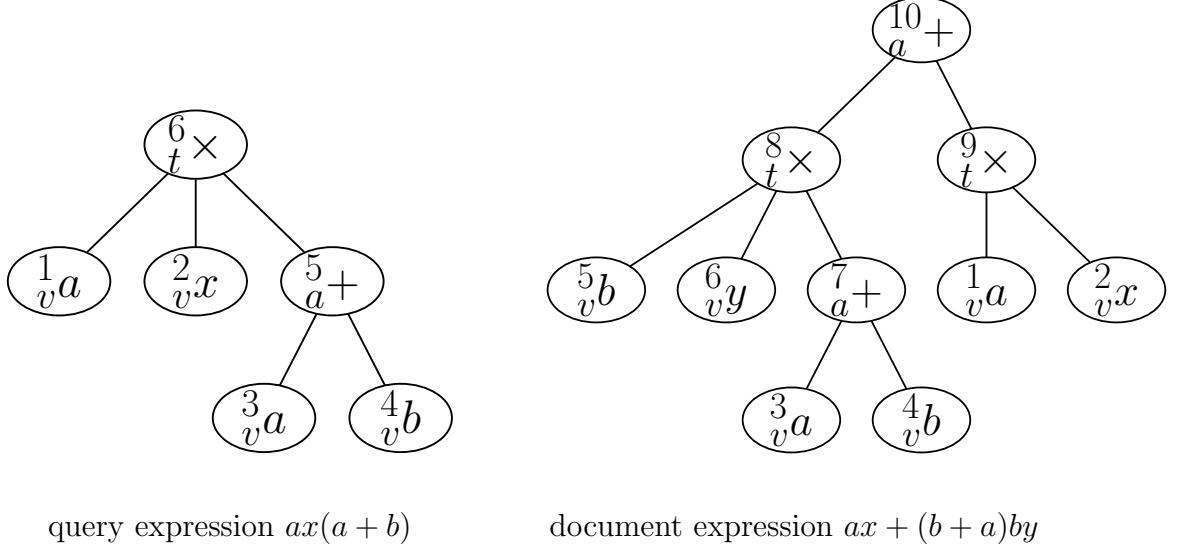
We may also need a way to evaluate how many percentage is a query matching a document expression. According to the property 2.3 in section 2.2.3 (observation #1), for two expression in formula subtree relation, we have  $\frac{|g(T_q)|}{|g(T_d)|} \leq 1$ , we name the ratio of left-hand as *matching-ratio*, which characterises the structural coverage for a matching query in an expression. Consider the scenario where a query expression is structural isomorphic to two document expressions, and the symbolic similarity between them are the same. However, the different “area” that the query matches to the two document expressions essentially makes them scored differently. For example, for query  $ax + b$ , document expression  $ax + b$  should precede  $x^2 + ax + b$  although the query matches both two document expressions equally in terms of symbol.

### 2.4.4 Final ranking schema

Before matching-depth and matching-ratio is introduced, if a formula subtree relation is inferred between two formula trees, they are considered structurally matching in a boolean manner. There is no similarity degree between a complete match (or a sub-formula) and structurally irrelevant. In another word, structural differences is considered in a way to filter indexed expressions, for further symbolic similarity assessment. The final output ranking is determined by symbolic similarity degree, without considering structural similarity. Matching-depth and matching-ratio are metres introduced to measure the structural similarity, in the end we need to combine them with symbolic similarity to rank search results. Denote symbolic similarity score calculated by MARK-AND-CROSS algorithm to be  $s$ , and  $f(d)$  being the matching depth factor,  $r$  being the matching ratio for a pair of query/document. We will use them together in a tuple  $(s, f(d), r)$ , to indicate the overall similarity and rank items in search results,

### 2.4.5 Illustrated by an example

We illustrate the method introduced in this chapter by a simple example here. Given a query expression  $ax(a + b)$  and document expression  $ax + (b + a)by$ , here we



**Figure 2.6:** Example query/document expression representation

show how we search the relevant document using this query and how the relevance score between them is calculated by mark-and-cross algorithm.

The query expression and document expression are represented by operation trees  $T_q$  and  $T_d$  in figure 2.6. Instead of only denoting the operation symbols at the internal nodes and the variable/constant symbols at the leaf nodes, we use the notation  $^iS$  or  $^iS$  to denote a node with symbol  $S$  labeled by  $l$  (i.e.  $\ell(S) = l$ ) with vertex number  $i$ . The three different possible labels used here are  $v$ ,  $t$  and  $a$ , standing for unified name “variable”, “times” and “add” respectively. To be concise and descriptive, we will interchangeably use either  $^iS$  notation or  $q_i$  (or  $p_i$ ) to represent the leaf-root path in a query (or document) operation tree where the leaf vertex  $i$  resides.

Firstly, the generated path sets for  $T_q$  and  $T_d$  are:

$$g(T_q) = \{1 \cdot 6, 2 \cdot 6, 3 \cdot 5 \cdot 6, 4 \cdot 5 \cdot 6\}$$

$$g(T_d) = \{5 \cdot 8 \cdot 10, 6 \cdot 8 \cdot 10, 3 \cdot 7 \cdot 8 \cdot 10, 4 \cdot 7 \cdot 8 \cdot 10, 1 \cdot 9 \cdot 10, 2 \cdot 9 \cdot 10\}$$

And the labeled path sets for each of the two are:

$$\begin{aligned}\ell(g(T_q)) &= \{v \cdot t, v \cdot a \cdot t\} \\ \ell(g(T_d)) &= \{v \cdot t \cdot a, v \cdot a \cdot t \cdot a\}\end{aligned}$$

Because  $\ell(g(T_q)) \cdot a = \ell(g(T_d))$ , which follows

$$\ell(g(T_q)) \cdot \hat{a} \subseteq \ell(g(T_d))$$

where  $\hat{a} = a$ , we know  $T_d$  is searchable by  $T_q$ , so there exists a path  $\hat{a}$  we can search to append after every labeled path in set  $\ell(g(T_d))$  so that we will find  $T_d$  by intersecting all the formula trees indexed (in index  $\Pi$ ) in these paths.

Secondly, by the implications from observation #1 of section 2.2.3, we get the candidate set for each of the path in  $T_q$ :

$$\begin{aligned}C_{q_1} &= \{p_5, p_6\} \\ C_{q_2} &= \{p_5, p_6\} \\ C_{q_3} &= \{p_3, p_4\} \\ C_{q_4} &= \{p_3, p_4\}\end{aligned}$$

In addition, get the list  $L$  containing all the query paths in  $T_q$  sorted by symbol and its occurrence in all path symbols,

$$\text{QList} = q_1, q_3, q_4, q_2.$$

Lastly, we can calculate the symbolic similarity degree between the two expressions by going through each query path from QList in order and apply mark-and-cross algorithm (use  $\alpha = 0.9$ ). Then the first two iterations which calculates the matching score for bound variable  $a$  in  $T_q$  can be illustrated by table 2.1.

Each path of query bound variable  $a$  is compared with that from document path in its candidate set, and each resulting  $\text{sim}(a, a')$  value is accumulated on the corresponding document bound variable  $B_v$ . Then the current score is also calculated, by accumulating the max  $B_v$  value among all the bound variable  $v$  in document expression. The tags associated with document path  $^4b$  and  $^5b$  are marked as crossed

| current score               | 0.9       |                   |         |             |
|-----------------------------|-----------|-------------------|---------|-------------|
| bound variable              | $B_a = 1$ | $B_b = 1.8$ (max) |         | $B_y = 0.9$ |
| document path<br>query path | ${}^3a$   | ${}^4b$           | ${}^5b$ | ${}^6y$     |
| ${}^1a$                     |           |                   | 0.9     | 0.9         |
| ${}^3a$                     | 1         | 0.9               |         |             |

**Table 2.1:** First two iterations of example score evaluation

| current score               | $0.9 + 0.9 = 1.8$ |           |         |           |
|-----------------------------|-------------------|-----------|---------|-----------|
| bound variable              | $B_a = 0.9$ (max) | $B_b = 0$ |         | $B_y = 0$ |
| document path<br>query path | ${}^3a$           | ${}^4b$   | ${}^5b$ | ${}^6y$   |
| ${}^1a$                     |                   |           | 0.9     | 0.9       |
| ${}^3a$                     | 1                 | 0.9       |         |           |
| ${}^4b$                     | 0.9               |           |         |           |

**Table 2.2:** 3rd iteration of example score evaluation



|                                           |                   |           |       |                   |
|-------------------------------------------|-------------------|-----------|-------|-------------------|
| <b>current score</b>                      | $1.8 + 0.9 = 2.7$ |           |       |                   |
| <b>bound variable</b>                     | $B_a = 0$         | $B_b = 0$ |       | $B_y = 0.9$ (max) |
| <b>document path</b><br><b>query path</b> | $^3a$             | $^4b$     | $^5b$ | $^6y$             |
| $^1a$                                     |                   |           | 0.9   | 0.9               |
| $^3a$                                     | 1                 | 0.9       |       |                   |
| $^4b$                                     | 0.9               |           |       |                   |
| $^2x$                                     |                   |           |       | 0.9               |

**Table 2.3:** 4th iteration of example score evaluation

state, to prevent path  $^4b$  and  $^5b$  from being compared in future iterations. For the next query bond variable  $b$  in QList, first zero every  $B_v$  value, then calculate its matching score with the document paths in its candidate sets in a similar manner except we have skipped some document paths in candidate sets because they are crossed in the previous iteration. Table 2.2 shows the intermediate results before iterate to the next query bond variable in QList. Finally, we use the same way to process the last query bond variable  $x$ . Its query path  $^2x$  is matched with document path  $^6y$  with  $\text{sim}(^2x, ^6y)$  value equals to 0.9. Then the only non-zero  $B_v$  value from document bound variable  $y$  will contribute to final symbolic similarity score which in the end is added up to 2.7 (see table 2.3).

Now we have finished our symbolic similarity evaluation between given query expression and document expression. We can also infer that the matching-depth  $d$  is  $|\hat{a}| = 1$ , matching factor  $f(d) = \frac{1}{1+d} = 0.5$ , and the matching-ratio  $r = \frac{|g(T_q)|}{|g(T_d)|} = \frac{4}{6} \approx 0.67$ . The final score tuple  $(2.7, 0.5, 0.67)$  is used to determine the similarity degree in general and rank the document expression.

## Chapter 3

### IMPLEMENTATION AND EVALUATION

Based on the method we propose, we have implemented a search engine as well as a crawler and a parser for  $\text{\LaTeX}$  content. Also, a demo <sup>1</sup> with Web interface is available for demonstration. In this paper, we will use this system to evaluate the effectiveness and efficiency of our system by comparing a complete system which measures symbolic similarity of expressions, with a baseline system which only does a boolean search of all structurally similar expressions for a query.

#### 3.1 Implementation

The outline of our system process is, crawl and parse  $\text{\LaTeX}$  content, and transform math expressions into an in-memory tree structure, write the leaf-root path labels along with degree number associated with each nodes (we call *branch word*) to disk as our index. The search engine search and score the similar expressions in index, and output the results. We have a search program which outputs results in standard out device, and a WEB front-end which outputs results in HTML pages.

##### 3.1.1 Crawler

##### 3.1.2 Parser

Because we are parsing  $\text{\LaTeX}$  directly, while most researches evaluate using MathML/XML format collection, and original  $\text{\LaTeX}$  information is not always preserved in the dataset converted from *LaTeXML*, we have tried to convert MathML/XML dataset back to  $\text{\LaTeX}$  using *pandoc*, but failed to convert all the document correctly.

---

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

| 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  | $\frac{e^x + y}{z}$                                   |
| 7  | $O(n \log n)$                                                                      | 8  | $H^n(X) = Z^n(X)/B^n(X)$                              |
| 9  | $A_n = \frac{1}{\pi} \int_{-\pi}^\pi F(x) \cos(nx) dx$                             | 10 | $\lim_{x \rightarrow \infty} (1 + \frac{1}{x})^x$     |
| 11 | $f(x) = f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \dots$                              | 12 | $f(a) = \frac{1}{2\pi i} \oint_r \frac{f(z)}{z-a} dz$ |
| 13 | $x^2 + 2xy + y^2 =  x ^2 + 2 x  y  +  y ^2$                                        | 14 | $\int_a^b f(x) dx = F(b) - F(a)$                      |
| 15 | $\frac{n!}{r_1! \cdot r_2! \cdot \dots \cdot r_k!}$                                | 16 | $-b \pm \sqrt{b^2 - 4ac}$                             |
| 17 | $1 + \tan^2 \theta = \sec^2 \theta$                                                | 18 | $\bar{u} = (x, y, z)$                                 |

**Table 3.1:** Test query set

Nevertheless, we have crawled *Mathjax* content (in  $\text{\LaTeX}$ ) from nearly all the posts from Math Stack Exchange website before March 2015 (over 60MB *bzip2* compressed raw data in text files with one  $\text{\LaTeX}$  formula per line) and choose to use our own dataset <sup>2</sup> to evaluate the data. Our test query set consists queries mostly from NTCIR-10 Math Formula Search [2] and [6], some of them we are not able to find similar formula in our own dataset, are excluded from our test query set. Table 3.1 shows our complete test queries used in our evaluation.

There are four levels of relevance scored from 0 to 4 in our evaluation. The

---

<sup>2</sup> data set can be downloaded from <http://infolab.ece.udel.edu:8912/cowpie/math-stack-exchange-march-2015.tar.bz2>

| 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           |

**Table 3.2:** Query statistics for matching structurally only

criteria considers both structural similarity and symbolic similarity. Structural similarity is measured by 0, 1 (mostly similar) and 2 (complete matching). While symbolic similarity is measured by 0, 1 (mostly similar for the matching parts) and 2 (identical symbols for the matching parts). The level of relevance is simply the sum of the two scores.

We have evaluated two methods in this paper, the first one is a boolean search for only structurally similar document expressions to the query. That is to say, every documents formula tree  $T_d \in \bigcap_{a \in L} \mathcal{I}_{\Pi}(a)$  is a hit, and only the first  $k$  hits are returned as search results where  $k$  is the maximum ranking items. Table 3.2 shows the distribution of hits and relevance level for each test query using this method. Another method we have evaluated, further applies mark-and-cross algorithm to score every formula tree in document expressions that structurally match the query. It uses a min-heap

| 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           |

**Table 3.3:** Query statistics for matching structurally with symbolic scoring

data structure that keeps replacing the minimal-score item in current ranking list with new one if the latter has a higher symbolic similarity score, which will give us top  $k$  highest-score hits with most symbolic similarity in our search result. Table 3.3 gives the distribution and relevance scores for this method. Both of the two method does not guarantee a thorough search through index (even if we are pruning and search only a set of the index), in fact, we set a limit (around 3,650,000) to the maximum items to be scanned in our index, to make our search response time practical for a real search engine.

## REFERENCES

- [1] Richard Zanibbi and Dorothea Blostein. Recognition and retrieval of mathematical expressions. *International Journal on Document Analysis and Recognition (IJDAR)*, 15(4):331–357, 2012.
- [2] 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>. Accessed: 2015-03-31.
- [3] Akiko Aizawa, Michael Kohlhase, and Iadh Ounis. Ntcir-11 math-2 task overview. *The 11th NTCIR Conference*, 2014.
- [4] Jozef Misutka. Mathematical search engine. Master’s thesis, Charles University in Prague, May 2013.
- [5] 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.
- [6] 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.
- [7] Qun Zhang and Abdou Youssef. An approach to math-similarity search. *Intelligent Computer Mathematics. International Conference, CICM*, 2014.
- [8] Miller B. and Youssef A. Technical aspects of the digital library of mathematical functions. *Annals of Mathematics and Artificial Intelligence* 38(1-3), 121136, 2003.
- [9] Youssef A. Information search and retrieval of mathematical contents: Issues and methods. *The ISCA 14th Intl Conf. on Intelligent and Adaptive Systems and Software Engineering (IASSE 2005)*, 2005.
- [10] Jozef Miutka and Leo Galambo. Extending full text search engine for mathematical content. *Towards Digital Mathematics Library.*, 2008.

- [11] Petr Sojka and Martin Lka. Indexing and searching mathematics in digital libraries. *Intelligent Computer Mathematics*, 6824:228–243, 2011.
- [12] Petr Sojka and Martin Lka. The art of mathematics retrieval. *ACM Conference on Document Engineering, DocEng 2011*, 2011.
- [13] Martin Lka. Evaluation of mathematics retrieval. Master’s thesis, Masarykova University, 2013.
- [14] P. Pavan Kumar, Arun Agarwal, and Chakravarthy Bhagvati. A structure based approach for mathematical expression retrieval. *Multi-disciplinary Trends in Artificial Intelligence*, 7694:23–34, 2012.
- [15] Robert Miner and Rajesh Munavalli. *An Approach to Mathematical Search Through Query Formulation and Data Normalization*. Springer Berlin Heidelberg, 2007.
- [16] Christopher D. Manning, Prabhakar Paghavan, and Hinrich Schutze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [17] Muhammad Adeel, Hui Siu Cheung, and Ar Hayat Khiyal. Math go! prototype of a content based mathematical formula search engine, 2008.
- [18] Andrea Asperti, Ferruccio Guidi, Claudio Sacerdoti Coen, Enrico Tassi, and Stefano Zacchiroli. A content based mathematical search engine: whelp. In *In: Post-proceedings of the Types 2004 International Conference, Vol. 3839 of LNCS*, pages 17–32. Springer-Verlag, 2004.
- [19] Michael Kohlhase and Ioan A. Sukan. A search engine for mathematical formulae. In *Proc. of Artificial Intelligence and Symbolic Computation, number 4120 in LNAI*, pages 241–253. Springer, 2006.
- [20] Michael Kohlhase. Mathwebsearch 0.4 a semantic search engine for mathematics.
- [21] Michael Kohlhase. Mathwebsearch 0.5: Scaling an open formula search engine.
- [22] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach (3rd Edition)*. Prentice Hall, December 2009.
- [23] Peter Graf. *Term Indexing*. Springer Verlag, 1996.
- [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] 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.



- [26] Hiroshi Ichikawa, Taiichi Hashimoto, Takenobu Tokunaga, and Hozumi Tanaka. New methods of retrieve sentences based on syntactic similarity. *IPSJ SIG Technical Reports, DBS-136, FI-79*, pages 39–46, 2005.
- [27] 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.
- [28] 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.
- [29] 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.
- [30] Shahab Kamali and Frank Wm. Tompa. A new mathematics retrieval system. *CIKM*, 2010.
- [31] 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.
- [32] 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.
- [33] 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.
- [34] Richard Zanibbi and Bo Yuan. Keyword and image-based retrieval for mathematical expressions. *Multi-disciplinary Trends in Artificial Intelligence. 6th International Workshop, MIWAI 2012.*, pages 23–34, 2011.
- [35] 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.
- [36] T. Nguyen, S. Hui, and K. Chang. A lattice-based approach for mathematical search using formal concept analysis. *Expert Systems with Applications*, 2012.
- [37] Akiko Aizawa, Michael Kohlhase, Iadh Ounis, and Moritz Schubotz. Ntcir-11 math-2 task overview. *Proc. of the 11th NTCIR Conference, Tokyo, Japan*, 2014.

- [38] Kamali Shahab and Tompa Frank Wm. Improving mathematics retrieval. *Towards a Digital Mathematics Library. Grand Bend, Ontario, Canada*, pages 37–48, 2009.
- [39] 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.
- [40] J. Roger Hindley. *Introduction to Combinators and [Lambda]-Calculus*. Cambridge University Press, 1986.

## **Appendix A**

### **TITLE OF APPENDIX A**

This is the information for the first appendix, Appendix A. Copy the base file, appA.tex, for each additional appendix needed such as appB.tex, appC.tex, etc. Modify the main base file to include each additional appendix file.

If there is only one appendix, then modify the main file to only use app.tex instead of appA.tex.

**Appendix B**  
**TITLE OF APPENDIX B**