

The Dyck Language Edit Distance Problem in Near-linear Time

Barna Saha

Dept. of Computer Science
University of Massachusetts, Amherst
Amherst, USA
barna@cs.umass.edu

Abstract—Given a string σ over alphabet Σ and a grammar G defined over the same alphabet, how many minimum number of repairs (insertions, deletions and substitutions) are required to map σ into a valid member of G ? The seminal work of Aho and Peterson in 1972 initiated the study of this *language edit distance problem* providing a dynamic programming algorithm for context free languages that runs in $O(|G|^2 n^3)$ time, where n is the string length and $|G|$ is the grammar size. While later improvements reduced the running time to $O(|G|n^3)$, the cubic time complexity on the input length held a major bottleneck for applying these algorithms to their multitude of applications.

In this paper, we study the language edit distance problem for a fundamental context free language, $\text{DYCK}(s)$ representing the language of well-balanced parentheses of s different types, that has been pivotal in the development of formal language theory. We provide the very first *near-linear time* algorithm to tightly approximate the $\text{DYCK}(s)$ language edit distance problem for any arbitrary s . $\text{DYCK}(s)$ language edit distance *significantly generalizes* the well-studied *string edit distance problem*, and appears in most applications of language edit distance ranging from data quality in databases, generating automated error-correcting parsers in compiler optimization to structure prediction problems in biological sequences. Its nondeterministic counterpart is known as the hardest context free language.

Our main result is an algorithm for edit distance computation to $\text{DYCK}(s)$ for any positive integer s that runs in $O(n^{1+\epsilon} \text{polylog}(n))$ time and achieves an approximation factor of $O(\frac{1}{\epsilon} \beta(n) \log |OPT|)$, for any $\epsilon > 0$. Here OPT is the optimal edit distance to $\text{DYCK}(s)$ and $\beta(n)$ is the best approximation factor known for the simpler problem of string edit distance running in analogous time. If we allow $O(n^{1+\epsilon} + |OPT|^2 n^\epsilon)$ time, then the approximation factor can be reduced to $O(\frac{1}{\epsilon} \log |OPT|)$. Since the best known near-linear time algorithm for the string edit distance problem has $\beta(n) = \text{polylog}(n)$, under near-linear time computation model both $\text{DYCK}(s)$ language and string edit distance problems have *polylog*(n) approximation factors. This comes as a surprise since the former is a significant generalization of the latter and their exact computations via dynamic programming show a stark difference in time complexity.

Rather less surprisingly, we show that the framework for efficiently approximating edit distance to $\text{DYCK}(s)$ can be utilized for many other languages. We illustrate this by considering various memory checking languages (studied extensively under distributed verification) such as STACK , QUEUE , PQ and DEQUE which comprise of valid transcripts of stacks, queues, priority queues and double-ended queues respectively. Therefore, any language that can be recognized by these data structures, can also be repaired efficiently by our algorithm.

Index Terms—edit distance; formal language; linear time algorithm design; approximation algorithms;

I. INTRODUCTION

Given a string σ over alphabet Σ and a grammar G defined over the same alphabet, how many minimum number of repairs (insertions, deletions and substitutions) are required to map σ into a valid member of G ? In this work, we consider such *language edit distance problem* with respect to $\text{DYCK}(s)$, where $|\Sigma| = 2s$. $\text{DYCK}(s)$ is a fundamental context free grammar representing the language of well-balanced parentheses of s different types, and DYCK language edit distance is a *significant generalization* of the *string edit distance problem* which has been studied extensively in theoretical computer science and beyond.

Dyck language appears in many contexts. These languages often describe a property that should be held by commands in most commonly used programming languages, as well as various subsets of commands/symbols used in LaTeX. Variety of semi-structured data from XML documents to JSON data interchange files to annotated linguistic corpora contain open and close tags that must be properly nested. They are frequently massive in size and exhibit complex structures with arbitrary levels of nesting tags (an XML document often encodes an entire database). For example, `dblp.xml` has current size of 1.2 GB, is growing rapidly, with 2.3 million articles that results in a string of parentheses of length more than 23 million till date. In addition, Dyck language plays an important role in DNA evolutionary languages and RNA structure modeling where the base nucleotide pairs in DNA/RNA sequences need to match up in a well-formed way. Deviations from this well-formed matching reveal interesting properties of the underlying biological sequences [17], [32]. Dyck language has been pivotal in the development of the theory of context-free languages (CFL). As stated by the Chomsky-Schutzenberger Theorem, every CFL can be mapped to a restricted subset of $\text{DYCK}(s)$ [10]. A comprehensive description of context free languages and Dyck languages can be found in [18], [21].

The study of *language edit distance problem* dates back to early seventies. Such an algorithm for context free grammar (CFG) was first proposed by Aho and Peterson that runs in $O(|G|^2 n^3)$ time where $|\sigma| = n$ is the string length and $|G|$ is the size of the grammar [1]. This was later improved by Myers to $O(|G|n^3)$ time [27]. These works were motivated by developing automated parsers for compiler design. For $\text{DYCK}(s)$

a dynamic programming algorithm gives a running time of $O(n^3)$ independent of s . Since a well-balanced string may be composed of two well-balanced substrings, the optimum edit distance computation for a substring from index i to j , $1 \leq i < j \leq n$, requires checking all possible decompositions of the substring at intermediate indices $k = i, i + 1, \dots, j$. This leads to the cubic dependency in the running time. It is possible to improve this bound slightly to $O(\frac{n^3}{\log n})$ using the Four-Russian technique [32].

Nearly two decades back, [26] reported these algorithms with cubic running time to be prohibitively slow for parser design. With modern data deluge, the issue of scalability has become far more critical. Motivated by a concrete application of repairing semi-structured documents where imbalanced parenthesis nesting is one of the major reported errors (14% of XML errors on the web is due to malformedness [16]) and lack of scalability of cubic time algorithms, the authors in [20] study the problem of approximating edit distance computation to $DYCK(s)$. Given any string σ if $\sigma' = \min_{x \in DYCK(s)} StrEdit(\sigma, x)$, they ask the question whether it is possible to design an algorithm that runs in near-linear time and returns $\sigma'' \in Dyck(s)$ such that $StrEdit(\sigma, \sigma'') \leq \alpha StrEdit(\sigma, \sigma')$ for some $\alpha \geq 1$ where $StrEdit$ is the normal string edit distance function and α is the approximation factor. Edit distance computation from a string of parentheses to $DYCK(s)$ is a significant generalization of string edit distance computation¹. A prototypical dynamic programming for string edit distance computation runs in quadratic time (as opposed to cubic time for $DYCK(s)$ edit distance problem). There is a large body of works on designing scalable algorithms for approximate string edit distance computation [4]–[7], [29]. Though basic in its appeal, nothing much is known for approximately computing edit distance to $DYCK(s)$.

In [20], the authors proposed fast greedy and branch and bound methods with various pruning strategies to approximately answer the edit distance computation to $DYCK(s)$, and showed its applicability in practice over rule based heuristics commonly employed by modern web browsers like Internet Explorer, Firefox etc. However, either their algorithms have worst-case approximation ratio as bad as $\Theta(n)$ or have running time exponential in $|OPT|$ (see [20] for worst case examples). It is to be noted that for $DYCK(1)$, there exists a simple single pass algorithm to compute edit distance: just pair up matching open and close parentheses and report the number of parentheses that could not be matched in this fashion.

In this paper, we study the question of approximating edit distance to $DYCK(s)$ for any $s \geq 2$ and give the *first* near-linear time algorithm with nontrivial approximation

¹For string edit distance computation, between string σ_1 and σ_2 over alphabet C , create a new alphabet $T \cup \bar{T}$ by uniquely mapping each character $c \in C$ to a new type of open parenthesis, say t_c , that now belongs to T . Let \bar{t}_c be the matching close parenthesis for t_c and we let $\bar{t}_c \in \bar{T}$. Now create strings σ'_1 by replacing each character of σ_1 with its corresponding open parenthesis in T , and create string σ'_2 by replacing each character of σ_2 with its corresponding close parenthesis in \bar{T} . Obtain σ by appending σ'_1 with reverse of σ'_2 . It is easy to check the edit distance between σ and $DYCK(s)$ is exactly equal to string edit distance between σ_1 and σ_2 .

guarantees. Our main result is an algorithm for edit distance computation to $DYCK(s)$ language for any positive integer s that runs in $O(n^{1+\epsilon} polylog(n))$ time and achieves an approximation factor of $O(\frac{1}{\epsilon} \beta(n) \log |OPT|)$, for any $\epsilon > 0$. Here OPT is the optimal edit distance to $DYCK(s)$ and $\beta(n)$ is the best approximation factor known for the simpler problem of string edit distance running in analogous time. If we allow $O(n^{1+\epsilon} + |OPT|^2 n^\epsilon)$ time, then the approximation factor can be reduced to $O(\log |OPT|)$. Since the best known near-linear time algorithm for the string edit distance problem has $\beta(n) = polylog(n)$, under near-linear time computation model both the string edit distance and the $DYCK(s)$ language edit distance problems have $polylog(n)$ approximation factors. This comes as a surprise since edit distance to $DYCK(s)$ is a significant generalization of the string edit distance problem and their exact computations via dynamic programming show a marked difference in time complexity.

Any parentheses string σ can be viewed as $Y_1 X_1 Y_2 X_2 \dots Y_z X_z$ for some $z \geq 0$ where Y_i s and X_i s respectively consist of only consecutive open and close parentheses. The special case of string edit distance problem can be cast as having $z = 1$. The approximation factor of our algorithm is in fact depends only on $\log z$, and not $\log |OPT|$. It is possible to ensure $z \leq OPT$ by a simple preprocessing. Our algorithm is based on a judicious combination of random walks in multiple phases that guide selection of subsequences having a single sequence of open parentheses followed by a single sequence of close parentheses. These subsequences are then repaired each by employing subroutine for STREDIT computation. The general framework of our algorithm and its analysis applies to languages far beyond $DYCK(s)$. We discuss this connection with respect to the memory checking languages whose study was initiated in the seminal work of Blum, Evans, Gemmell, Kannan and Naor [8] with numerous follow-up works [2], [9], [11], [13], [28]. We consider basic languages such as STACK, QUEUE, PQ, DEQUE etc. They comprise of valid transcripts of stacks, queues, priority queues and double-ended queues respectively. Given a transcript of any such memory-checking language, we consider the problem of finding the minimum number of edits required to make the transcript error-free and show that the algorithm for $DYCK(s)$ can be adapted to return a valid transcript efficiently with the same approximation bound. Therefore, any language that can be recognized by these data structures, can also be repaired efficiently by our algorithm. We believe our novel multi-phase random walk based technique can be useful for any generic sequence alignment type problems and may lead to a systematic way of speeding up dynamic programming algorithms for this large class of problems.

A. Related Work

Early works on edit distance to grammar [1], [27] was motivated by the problem of correcting and recovering from syntax error during context-free parsing and have received significant attention in the realm of compiler optimization [14], [15], [19], [26]. Many of these works focus on designing time-

efficient parsers using local recovery [15], [19], [26] rather than global dynamic programming based algorithms [1], [27], but to the best of our knowledge, none of these methods provide approximation guarantee on edit distance in sub-cubic time. Approximating edit distance to $\text{DYCK}(s)$ has recently been studied in [20] for repairing XML documents, and in [32] for RNA folding. But either the proposed algorithms have running time very close to cubic [32], or the proposed subcubic algorithms all have worst case approximation factor $\Theta(n)$ [20].

Recognizing a grammar is a much simpler task than repairing, and DYCK language recognition has attracted considerable attention before. Using a stack, it is straightforward to recognize $\text{DYCK}(s)$ in a single pass—a prototypical example of a stack-based algorithm. When there is a space restriction, Magniez, Mathieu and Nayak [25] considered the streaming complexity of recognizing $\text{DYCK}(s)$ showing an $\Omega(\sqrt{n})$ lower bound and a matching upper bound within a $\log n$ factor. Even with multiple one-directional passes, the lower bound remains at $\Omega(\sqrt{n})$ [9], surprisingly with two passes in opposite directions the space complexity reduces to $O(\log^2 n)$. This exhibits a curious phenomenon of streaming complexity of recognizing $\text{DYCK}(s)$. Krebs et al. extended the work of [25] to consider very restricted cases of errors, where an open (close) parenthesis can only be changed into another open (close) parenthesis [22]. Again with only such edits, computing edit distance in linear time is trivial: whenever an open parenthesis at the stack top cannot be matched with the current close parenthesis, change one of them to match. Allowing arbitrary insertions, deletions and substitutions is what makes the problem significantly harder. In the property testing framework, in the seminal paper Alon, Krivelevich, Newman and Szegedy [3] showed that $\text{DYCK}(1)$ is testable in time independent of n , however $\text{DYCK}(2)$ requires $\Omega(\log n)$ queries. This lower bound was further strengthened to $n^{1/11}$ by Parnas, Ron and Rubinfeld [30] where they also give a testing algorithm using $n^{2/3}/\epsilon^3$ queries. These algorithms can only distinguish between the case of 0 error with ϵn errors, and therefore, are not applicable to the problem of approximating edit-distance to $\text{DYCK}(s)$.

Edit distance to $\text{DYCK}(s)$ is a significant generalization of string edit distance problem. String edit distance enjoys the special property that if the optimal edit distance is d then a symbol at the i th position in one string must be matched with a symbol at a position between $(i - d)$ to $(i + d)$ in the other string, if d is the optimum string edit distance. Using this “local” property, prototypical quadratic dynamic programming algorithm can be improved to run in time $O(dn)$ which was later improved to $O(n+d^7)$ [31] to $O(n+d^4)$ [12] and then to $O(n + d^2)$ [23]. The later result implies a \sqrt{n} -approximation for string edit distance problem. However, all of these crucially use the locality property, which does not hold for parenthesis strings: two parentheses far apart can match as well. Also, it is known that parsing arbitrary CFG is as hard as boolean matrix multiplication [24] and a nondeterministic version of DYCK is the hardest CFG [10]. Therefore, exactly computing edit distance to $\text{DYCK}(s)$ in time much less than subcubic

would be a significant accomplishment. For string edit distance, the current best approximation ratio of $(\log n)^{O(\frac{1}{2})}$ in $O(n^{1+\epsilon} \text{polylog}(n))$ running time for any fixed $\epsilon > 0$ is due to Andoni, Krauthgamer and Onak [4]. This result is preceded by a series of works which improved the approximation ratio from \sqrt{n} [23] to $n^{\frac{3}{7}}$ [6], then to $n^{\frac{1}{3}+o(1)}$ [7], all of which run in linear time and finally to $2^{\sqrt{\log n \log \log n}}$ that run in time $n2^{\sqrt{\log n \log \log n}}$ [5].

B. Techniques & RoadMap

Definition 1. The congruent of a parenthesis x is defined as its symmetric opposite parenthesis, denoted \bar{x} . The congruent of a set of parentheses X , denoted \bar{X} , is defined as $\{\bar{x} \mid x \in X\}$.

We use T to denote the set of open parentheses and \bar{T} to denote the set of close parentheses. (Each $x \in T$ has exactly one congruent $\bar{x} \in \bar{T}$ that it matches to and vice versa.) The alphabet $\Sigma = T \cup \bar{T}$.

Definition 2. A string over some parenthesis alphabet $\Sigma = T \cup \bar{T}$ is called well-balanced if it obeys the context-free grammar $\text{DYCK}(s)$, $s = |T|$ with productions $S \rightarrow SS$, $S \rightarrow \phi$ (empty string) and $S \rightarrow aS\bar{a}$ for all $a \in T$.

Definition 3. The DYCK Language Edit Distance Problem, given string $\sigma = \sigma_1 \dots \sigma_n$ over alphabet $\Sigma = T \cup \bar{T}$, is to find $\arg \min_{\sigma'} \text{StrEdit}(\sigma, \sigma')$ such that $\sigma' \in \text{DYCK}(s)$, $s = |T|$.

A Simple Algorithm (Section II): We start with a very simple algorithm that acts as a stepping stone for the subsequent refinements. The algorithm is as simple as it gets, and is referred to as **Random-deletion**.

*Initialize a stack to empty. Scan the parenthesis string σ left to right. If the current symbol is an open parenthesis, insert it into the stack. If the current symbol is a close parenthesis, check the symbol at top of the stack. If both the symbols can be matched, match them. If the stack is empty, delete the current symbol. Else **delete one of them independently with equal probability $\frac{1}{2}$** . If the stack is nonempty when the scan has ended, delete all symbols from the stack.*

Let d be the optimum edit distance to $\text{DYCK}(s)$. Each deletion accounts for one edit. We show

Theorem 1. Random-deletion obtains a $4d$ -approximation for edit distance computation to $\text{DYCK}(s)$ for any $s \geq 2$ in linear time with constant probability.

The probability can be boosted by running the algorithm $\Theta(\log n)$ times and considering the iteration which results in the minimum number of edits. In the worst case, when $d = \sqrt{n}$, the approximation factor can be $4\sqrt{n}$. This also gives a very simple algorithm for string edit distance problem that achieves a $O(\sqrt{n})$ approximation.

The analysis of even such a simple algorithm is nontrivial and proceeds as follows. First, we allow the optimal algorithm to consider only deletion and allow it to process the string using a stack; this increases the edit distance by a factor of 2 (Lemma 6, Lemma 7). We then define for each comparison where an open and close parenthesis cannot be matched, a

corresponding correct and wrong move. If an optimal algorithm also compares exactly the two symbols and decides to delete one, then we can simply define the parenthesis that is deleted by the optimal algorithm as the correct move and the other as a wrong move. However, after the “first” wrong move, the comparisons performed by our procedure vs the optimal algorithm can become very different. Yet, we can label one of the two possible deletions as a correct move in some sense of decreasing distance to an optimal state. Consider a toy example, $cb\bar{a}b\bar{c}$. An optimal algorithm deletes \bar{a} to make the string well-formed. The first comparison performed by our algorithm may take a wrong move by deleting b instead of \bar{a} . While comparing c and \bar{a} , it still can recover from this earlier error without a big penalty on the edit cost if it deletes \bar{a} -the correct move. We show that if up to time t , the algorithm has taken W_t wrong moves then it is possible to get a well-formed string using a total of $2(d + 2W_t)$ edit operations. These two properties help us to map the process of deletions to a one dimensional random walk-the GAMBLER’S RUIN problem.

In the considered gambler’s ruin problem, a gambler enters a casino with $\$d$ (remember d is the optimum edit distance to $\text{DYCK}(s)$) and starts playing a game where he wins with probability $1/2$ and loses with probability $1/2$ independently. The gambler plays the game repeatedly betting $\$1$ in each round. He leaves if he runs out of money or gets $\$n$. We can show that the number of edit operations performed by our algorithm can not be more than the number of steps taken by the gambler to be ruined. However, on expectation, the gambler takes $O(n)$ steps to be ruined, and this bound is not useful for our purpose. Interestingly, the underlying probability distribution of the number of steps taken by the gambler is heavy-tailed and using that property, one can still show that there is considerable probability ($\sim \frac{1}{5}$) that gambler is ruined in $O(d^2)$ steps. Therefore the total number of edits of our algorithm is bounded by $O(d^2)$ leading to an $O(d)$ approximation.

A Refined Algorithm (Section III): We now refine our algorithm as follows. Given string σ , we can delete any prefix of close parentheses, delete any suffix of open parentheses and match well-formed substrings without affecting the optimal solution. After that, σ can be written as $Y_1X_1Y_2X_2\dots Y_zX_z$ where each Y_i is a sequence of open parentheses, each X_i is a sequence of close parentheses and $z \leq d$. In the optimal solution X_1 is matched with some suffix of Y_1 possibly after doing edits. Let us denote this suffix by Z_1 . If we can find the left boundary of Z_1 , then we can employ STREDIT to compute string edit distance between Z_1 and X_1 (we need to consider reverse of X_1 and convert each $t \in X_1$ to \bar{t} -this is what is meant when we refer STREDIT between a sequence of open and a sequence of close parentheses), as Z_1X_1 consists of only a single sequence of open parentheses followed by a single sequence of close parentheses. If we can identify Z_1 correctly, then in the optimal solution X_2 is matched with a suffix of $Y_1^{res}Y_2$ where $Y_1^{res} = Y_1 \setminus Z_1$. Let us denote it by Z_2 . If we can again find the left boundary of Z_2 , then we can employ STREDIT between Z_2 and X_2 and so on. The question

is *how do we compute these boundaries?*

The key trick is to use *Random-deletion* again (repeat it appropriately $\sim \log n$ times) to find these boundaries approximately (see Section III). We consider the suffix of Y_1 which *Random-deletion* matches against X_1 possibly after multiple deletions (call it Z'_1) and use Z'_1 to approximate Z_1 . We show again using the mapping of *Random-deletion* to random walk, that the error in estimating the left boundary is bounded (Lemma 9, Lemma 16, Lemma 17). Specifically, if $\text{StrEdit}(Z_1, X_1) = d_1$, then the error in estimating the boundaries is at most $d_1\sqrt{2\log d_1}$ and $\text{StrEdit}(Z'_1, X_1) \leq 2d_1\sqrt{2\ln d_1}$. But, the error that we make in estimating Z_1 may propagate and affect the estimation of Z_2 . Hence the gap between optimal Z_2 and estimated Z'_2 becomes wider. If $\text{StrEdit}(Z_2, X_2) = d_2$, then we get $\text{StrEdit}(Z'_2, X_2) \leq 2(d_1 + d_2)\sqrt{2\ln(d_1 + d_2)}$. Proceeding, in this fashion, we get the following theorem.

Theorem 2. *The refined algorithm obtains an $O(z\beta(n)\sqrt{\ln d})$ -approximation factor for edit distance computation from strings to $\text{DYCK}(s)$ for any $s \geq 2$ in $O(n \log n + \alpha(n))$ time with probability at least $(1 - \frac{1}{n} - \frac{1}{d})$, where there exists an algorithm for STREDIT running in $\alpha(n)$ time that achieves an approximation factor of $\beta(n)$.*

Further Refinement: Main Algorithm (Section IV): Is it possible to compute subsequences of σ such that each subsequence contains a single sequence of open and close parentheses in order to apply STREDIT, yet propagational error can be avoided? This leads to our main algorithm.

Example. Consider $\sigma = Y_1X_1Y_2X_2Y_3X_3\dots Y_zX_z$, and let the optimal algorithm matches X_1 with $Z_{1,1}$, matches X_2 with $Z_{1,2}Y_2$, matches X_3 with $Z_{1,3}Y_3$, and so on, where $Y_1 = Z_{1,z}Z_{1,z-1}\dots Z_{1,2}Z_{1,1}$. In the *refined algorithm*, when *Random-deletion* finishes processing X_1 and tries to estimate the left boundary of $Z_{1,1}$, it might have already deleted some symbols of $Z_{1,2}$. It is possible that it deletes $\Theta(d_1\sqrt{\log d_1})$ symbols from $Z_{1,2}$. Therefore, when computing STREDIT between X_1 and Z'_1 , the portion of $Z_{1,2}$ in Z'_1 may not have any matching symbols and results in an increased STREDIT computation. More severely, the *Random-deletion* process gets affected when processing X_2 . *Random-deletion* does not find the already deleted symbols in $Z_{1,2}$ which ought to be matched with some subsequence of X_2 . As a result, it may start comparing $Z_{1,3}$ with X_2 , and in the process may delete up to $\Theta((d_1 + d_2)\sqrt{\log(d_1 + d_2)})$ symbols from $Z_{1,3}$, and so on. To remedy this, view $X_2 = X_{2,in}X_{2,out}$ where $X_{2,in}$ is the prefix of X_2 that is matched with Y_2 and $X_{2,out}$ is matched with $Z_{1,2}$. Consider pausing the random deletion process when it finishes Y_2 and thus attempt to find $X_{2,in}$. Suppose, *Random-deletion* matches $X'_{2,in}$ with Y_2 , then compute $\text{StrEdit}(Y_2, X'_{2,in})$. While there could still be a mistake in computing $X'_{2,in}$, the mistake does not affect $Z_{1,3}$. Else if *Random-deletion* process finishes X_2 before finishing Y_2 , then of course it affected $Z_{1,3}$. In that case Z'_2 is a suffix of Y_2 and we compute $\text{StrEdit}(Z'_2, X_2)$. Similarly, when processing X_3 , we pause whenever X_3 or Y_3 is exhausted and

create an instance of STREDIT accordingly. Suppose, for the sake of this example, Y_2, Y_3, \dots, Y_z are finished before finishing X_2, X_3, \dots, X_z respectively and X_1 is finished before Y_1 . Then, after creating the instances of STREDITS as described, we are left with a sequence of open parenthesis corresponding to a prefix of Y_1 and a sequence of close parenthesis, which is a combination of suffixes from X_2, X_3, \dots, X_z . We can compute STREDIT between them. Since most of $Z_{1,z}Z_{1,z-1}\dots Z_{1,2}$ exists in this remaining prefix of Y_1 and their matching parentheses in the constructed sequence of close parentheses, the computed STREDIT distance remains small.

Let us call each X_iY_i a block. As the example illustrates, we create STREDIT instances corresponding to what *Random-deletion* does *locally* inside each block. After the first phase, from each block we are either left with a sequence of open parentheses (call it a block of type O), or a sequence of close parentheses (call it a block of type C), or the block is empty. This creates new sequences of open and close parentheses by combining all the consecutive O blocks together (remove empty blocks) and similarly combining all consecutive C blocks together (remove empty blocks). We get at most $\lfloor \frac{z}{2} \rfloor$ new blocks in the remaining string after deleting any prefix of close and any suffix of open parentheses. We now repeat the process on this new string. This process can continue at most $\lceil \log z \rceil + 1$ phases, since the number of blocks reduces at least by a factor of 2 going from one phase to the next. This entire process is repeated $O(\log n)$ time and the final outcome is the minimum of the edit distances computed over these repetitions. The following theorem summarizes the performance of this algorithm.

Theorem 3. *There exists an algorithm that obtains an $O(\beta(n) \log z \sqrt{\ln d})$ -approximation factor for edit distance computation to $\text{DYCK}(s)$ for any $s \geq 2$ in $O(n \log n + \alpha(n))$ time with probability at least $(1 - \frac{1}{n} - \frac{1}{d})$, where there exists an algorithm for STREDIT running in $\alpha(n)$ time that achieves an approximation factor of $\beta(n)$, and z is the number of blocks.*

The $\sqrt{\log d}$ factor in the approximation can be avoided if we consider iterating $O(n^\epsilon \log n)$ times. Since, the best known near-linear time algorithm for STREDIT anyway has $\alpha(n) = n^{1+\epsilon}$ and $\beta = (\log n)^{\frac{1}{2}}$, we obtain the following theorem.

Theorem 4. *For any $\epsilon > 0$, there exists an algorithm that obtains an $O(\frac{1}{\epsilon} \log z (\log n)^{\frac{1}{2}})$ -approximation factor for edit distance computation to $\text{DYCK}(s)$ for any $s \geq 2$ in $O(n^{1+\epsilon})$ time with high probability, and z is the number of blocks.*

If instead we apply the string edit distance computation algorithm [23] in Theorem 3 and Theorem 4, we get the corollary

Corollary 5. (i) *There exists an algorithm that obtains an $O(\log z \sqrt{\ln d})$ -approximation factor for edit distance computation to $\text{DYCK}(s)$ for any $s \geq 2$ in $O(n + d^2)$ time with high probability, and also compute the edits.*

(ii) *There exists an algorithm that obtains an $O(\frac{1}{\epsilon} \log z)$ -*

approximation factor for edit distance computation to $\text{DYCK}(s)$ for any $s \geq 2$ in $O(n^{1+\epsilon} + d^2 n^\epsilon)$ time with high probability, and also compute the edits.

The algorithm and its analysis gives a general framework which can be applied to many other languages beyond $\text{DYCK}(s)$. Employing this algorithm one can repair footprints of several memory checking languages such as STACK, QUEUE, PQ and DEQUE efficiently. This extension can be found in the full version.

II. ANALYSIS OF **Random-deletion**

Here we analyse the performance of **Random-deletion** and prove Theorem 1.

We consider only deletion as a viable edit operation and under deletion-only model, assume that the optimal algorithm is stack based, and matches well-formed substrings greedily. The following two lemmas state that we lose only a factor 2 in the approximation by doing so. Note that all the missing proofs can be found in the full version.

Lemma 6. *For any string $\sigma \in (T \cup \bar{T})^*$, $\text{OPT}(\sigma) \leq \text{OPT}_d(\sigma) \leq 2\text{OPT}(\sigma)$, where $\text{OPT}(\sigma)$ is the minimum number of edits: insertions, deletions, substitutions required and $\text{OPT}_d(\sigma)$ is the minimum number of deletions required to make σ well-formed.*

Lemma 7. *There exists an optimal algorithm that makes a single scan over the input pushing open parentheses to stack and on observing a close parenthesis, the algorithm compares it with the stack top. If the symbols match, then both are removed from further consideration, otherwise one of the two symbols is deleted.*

From now onward we fix a specific optimal stack based algorithm, and refer that as the optimal algorithm.

Let us initiate time $t = 0$. At every step in *Random-deletion* when we either match two parentheses (current close parenthesis in the string with open parenthesis at the stack top) or delete one of them, we increment time t by 1.

We define two sets A_t and A_t^{OPT} for each time t .

Definition 4. *For every time $t \geq 0$, A_t is defined as all the indices of the symbols that are matched or deleted by **Random-deletion** up to and including time t .*

Definition 5. *For every time $t \geq 0$, $A_t^{OPT} = \{i \mid i \in A_t \text{ or } i \text{ is matched by the optimal algorithm with some symbol with index in } A_t\}$.*

Clearly at all time $t \geq 0$, $A_t^{OPT} \supseteq A_t$. We now define a correct and wrong move.

Definition 6. *A comparison at time t in the algorithm leading to a deletion is a correct move if $|A_t^{OPT} \setminus A_t| \leq |A_{t-1}^{OPT} \setminus A_{t-1}|$ and is a wrong move if $|A_t^{OPT} \setminus A_t| > |A_{t-1}^{OPT} \setminus A_{t-1}|$.*

Lemma 8. *At any time t , there is always a correct move, and hence **Random-deletion** always takes a correct move with probability at least $\frac{1}{2}$.*

Proofsketch. Suppose the algorithm compares an open parenthesis $\sigma[i]$ with a close parenthesis $\sigma[j]$, $i < j$ at time t , and they do not match. If possible, suppose that there is no correct move. Since *Random-deletion* is stack-based, A_t contains all indices in $[i, j]$. It may also contain intervals of indices $[i_1, j_1], [i_2, j_2], \dots$ because there can be multiple blocks. It must hold $[1, j-1] \setminus A_t$ does not contain any close parenthesis. Now for both the two possible deletions to be wrong, the optimal algorithm must match $\sigma[i]$ with some $\sigma[j']$, $j' > j$, and also match $\sigma[j]$ with $\sigma[i']$, $i' < i, i' \notin A_t$. But, this is not possible due to the property of well-formedness. Similarly, it can be argued that a correct move exists if at time t either the stack is empty or the string is exhausted. \square

Lemma 9. *If at time t (up to and including time t), the number of indices in A_t^{OPT} that the optimal algorithm deletes is d_t and the number of correct and wrong moves are respectively c_t and w_t then $|A_t^{OPT} \setminus A_t| \leq d_t + w_t - c_t$.*

The proof considers various possible states of the algorithm at time t , and compares A_t^{OPT} with A_t to obtain the bound. Let S_t denote the string σ at time t after removing all the symbols that were deleted by **Random-deletion** up to and including time t .

Lemma 10. *Consider d to be the optimal edit distance to DYCK(s). If at time t (up to and including t), the number of indices in A_t^{OPT} that the optimal algorithm deletes be d_t and $|A_t^{OPT} \setminus A_t| = r_t$, at most $r_t + (d - d_t)$ edits is sufficient to convert S_t into a well-balanced string.*

Proof. Since $|A_t^{OPT} \setminus A_t| = r_t$, there exists exactly r_t indices in S_t such that if those indices are deleted, the resultant string is same as what the optimal algorithm obtains after processing the symbols in A_t^{OPT} . For the symbols in remaining $\{1, 2, \dots, n\} \setminus A_t^{OPT}$, the optimal algorithm does at most $d - d_t$ edits. Therefore a total of $r_t + (d - d_t)$ edits is sufficient to convert S_t into a well-balanced string. \square

Lemma 11. *The edit distance between the final string S_n and σ is at most $d + 2w_n$.*

Proof. Consider any time $t \geq 0$, if at t , the number of deletions by the optimal algorithm in A_t^{OPT} is d_t , the number of correct moves and wrong moves are respectively c_t and w_t , then we have $|A_t^{OPT} \setminus A_t| \leq d_t + w_t - c_t$. The number of edits that have been performed to get S_t from S_0 is $c_t + w_t$. Denote this by $E(0, t)$. The number of edits that are required to transform S_t to well-formed is at most $(d - d_t) + d_t + w_t - c_t = d + w_t - c_t$ (by Lemma 10). Denote it by $E'(t, n)$. Hence the minimum total number of edits required (including those already performed) considering state at time t is $E(0, t) + E'(t, n) = d + 2w_t$. Since this holds for all time t , the lemma is established. \square

In order to bound the edit distance, we need a bound on w_n . To do so we map the process of deletions by **Random-deletion** to a random walk.

A. Mapping into Random Walk

We consider the following one dimensional random walk. The random walk starts at coordinate d , at each step, it moves one step right (+1) with probability $\frac{1}{2}$ and moves one step left (-1) with probability $\frac{1}{2}$. We count the number of steps required by the random walk to hit the origin.

We now associate a modified random walk with the deletions performed by *Random-deletion* as follows. Every time *Random-deletion* needs to take a move (performs one deletion), we consider one step of the modified random walk. If *Random-deletion* takes a wrong move, we let this random walk make a right (away from origin) step. On the other hand if *Random-deletion* takes a correct move, we let this random walk take a left step (towards origin move). If the random walk takes W right steps, then *Random-deletion* also makes W wrong moves. If the random walk takes W right steps before hitting the origin, then it takes in total a $d + 2W$ steps, and *Random-deletion* also deletes $d + 2W$ times. Therefore, hitting time of this modified random walk starting from d characterizes the number of edit operations performed by *Random-deletion*. In this random walk, left steps (towards origin) are taken with probability $\geq \frac{1}{2}$ (sometimes with probability 1). Therefore, hitting time of this modified random walk is always less than the hitting time of an one-dimensional random walk starting at d and taking right and left step independently with equal probability.

We therefore calculate the probability of a one-dimensional random walk taking right or left steps with equal probability to have a hitting time D starting from d . The computed probability serves as a lower bound on the probability that *Random-deletion* takes D edit operations to transform σ to well-formed. This one dimensional random walk is known as GAMBLER'S RUIN problem. We are interested in the probability that gambler gets ruined.

We now calculate the probability that the gambler is ruined in D steps precisely. Let \mathcal{P}_d denote the law of a random walk starting in $d \geq 0$, let $\{Y_i\}_0^\infty$ be the i.i.d. steps of the random walk, let $S_D = d + Y_1 + Y_2 + \dots + Y_D$ be the position of random walk starting in position d after D steps, and let $T_0 = \inf D : S_D = 0$ denotes the walks first hitting time of the origin. Clearly $T_0 = 0$ for \mathcal{P}_0 . Then we can show

Lemma 12. *For the GAMBLER'S RUIN problem $\mathcal{P}_d(T_0 = D) = \frac{d}{D} \left(\frac{D-d}{2} \right) \frac{1}{2^D}$.*

We now calculate the probability that a gambler starting with $\$d$ hits 0 within cd steps. Our goal will be to minimize c as much as possible, yet achieving a significant probability of hitting 0.

Lemma 13. *In GAMBLER'S RUIN, the gambler starting with $\$d$ hits 0 within $2d^2$ steps with probability at least 0.194.*

Corollary 14. *In GAMBLER'S RUIN, the gambler starting with $\$d$ hits 0 within $\frac{1}{\epsilon} \frac{d^2}{\log d}$ steps for any constant $\epsilon > 0$ with probability at least $\frac{\sqrt{\epsilon \log d}}{d^\epsilon}$.*

Theorem [1] follows from the mapping that the edit distance computed by *Random-deletion* is at most the number of steps taken by gambler's ruin to hit the origin starting from $\$d$ and then applying Lemma 13 and noting that we are only using deletions (Lemma 6 and Lemma 7).

III. ANALYSIS OF THE REFINED ALGORITHM

The algorithm proceeds as follows. It continues **Random-deletion** process as before, but now it keeps track of the substring with which each X_a , $a = 1, 2, \dots, z$ is matched (possibly through multiple deletions) during this random deletion process. While processing X_1 , the random deletion process is restarted $3 \log_b n$ times, $b = \frac{1}{(1-0.194)} = 1.24$ and at each time the algorithm keeps a count on how many deletions are necessary to complete processing of X_1 . It then selects the particular iteration in which the number of deletions is minimum. We let $Z_{1,min}$ to denote the substring to which X_1 is matched in that iteration. The algorithm then continues the random deletion process. It next stops when processing on X_2 finishes. Again, the portion of random deletion process between completion of processing X_1 and completion of processing X_2 is repeated $3 \log_b n$ times and the iteration that results in minimum number of deletions is picked. We define $Z_{2,min}$ accordingly. The algorithm keeps proceeding in a similar manner until the string is exhausted. In the process, $Z_{a,min}$ is matched with X_a for $a = 1, 2, \dots, z$. But, instead of using the edits that the random deletion process makes to match $Z_{a,min}$ to X_a , our algorithm invokes the best string edit distance algorithm $StrEdit(Z_{a,min}, X_a)$ which converts $Z_{a,min}$ to R_a and X_a to T_a such that $R_a T_a$ is well-formed. Clearly, at the end we have a well-formed string.

A. Analysis

We first analyze its running time.

Lemma 15. *The expected running time of the refined algorithm is $O(n \log n + \alpha(n))$ where $\alpha(n)$ is the running time of $StrEdit$ to approximate string edit distance of input string of length n within factor $\beta(n)$.*

We now proceed to analyze the approximation factor. For that we assume that the optimal distance d is at least 3 (otherwise employ the algorithm of the previous section). Let the optimal edit distance to convert $Z_a X_a$ into well-formed be d_a for $a = 1, 2, \dots, z$ where $Z_a = Z_{1,a} Z_{2,a} \dots Z_{a,a}$.

While computing the set $Z_{a,min}$, it is possible that our algorithm inserts symbols outside of Z_a to it or leaves out some symbols of Z_a . In the former case, among the extra symbols that are added, if the optimal algorithm deletes some of these symbols as part of some other $Z_{a'}, a' \neq a$, then these deletions are "harmless". If we only include these extra symbols to $Z_{a,min}$, then we can as well pretend that those symbols are included in Z_a too. The edit distance of the remaining substrings are not affected by this modification. Therefore, for analysis purpose, **both for this algorithm and for the main algorithm in the next section, we always**

assume w.l.o.g that the optimal algorithm does not delete any of the extra symbols that are added.

Lemma 16. *The number of deletions made by random deletion process to finish processing X_1, X_2, \dots, X_l , for $l = 1, 2, \dots, z$, that is to match $Z_{a,min}, X_a$, $a \leq l$, is at most $2(\sum_{a=1}^l d_a)^2$ with probability at least $(1 - \frac{1}{n^3})^l$.*

Let us denote by C_l and W_l the number of correct and wrong moves taken by the random deletion process when processing of X_1, X_2, \dots, X_l finishes. Since at each deletion, correct move has been taken with probability at least $\frac{1}{2}$ then by the standard Chernoff bound argument followed by union bound we have the following lemma.

Lemma 17. *When the processing of X_1, X_2, \dots, X_l finishes $W_l - C_l \leq (2 \sum_{a=1}^l d_a) \sqrt{2 \ln d}$ with probability at least $(1 - \frac{1}{n} - \frac{1}{d})$.*

Proof. Probability that the number of deletions made by *Random-deletion* process is at most $2(\sum_{a=1}^l d_a)^2$ is $\geq (1 - \frac{1}{n^3})^l$. Let us denote the number of these deletions by D_l , for $l = 1, 2, \dots, z$. Hence $D_l = W_l + C_l$. We use the following version of Azuma's inequality for simple coin flips to bound $W_l - C_l$.

Azuma's inequality for coin flips. Let F_i be a sequence of independent and identically distributed random coin flips taking values -1 or 1 . Defining $X_i = \sum_{j=1}^i F_j$ yields a martingale with $|X_k - X_{k+1}| \leq 1$, and Azuma's inequality states $\Pr[X_N > t] \leq \exp(-\frac{t^2}{2N})$.

We are only interested to bound $W_l - C_l$ (and not the absolute difference), and the worst case bound occurs when wrong and correct moves are taken with equal probability. Hence $\Pr[W_l - C_l > (2 \sum_{a=1}^l d_a) \sqrt{2 \ln d} | D_l \leq 2d_l^2]$

is at most $\exp(-\frac{8(\sum_{a=1}^l d_a)^2 \ln d}{4(\sum_{a=1}^l d_a)^2}) = \frac{1}{d^2}$. We have $\Pr[W_l - C_l > (2 \sum_{a=1}^l d_a) \sqrt{2 \ln d}] \leq \Pr[D_l > 2d_l^2] + \Pr[W_l - C_l > (2 \sum_{a=1}^l d_a) \sqrt{2 \ln d} | D_l \leq 2d_l^2] \Pr[D_l \leq 2d_l^2] \leq 1 - (1 - \frac{1}{n^3})^l + (1 - \frac{1}{n^3})^l \frac{1}{d^2} \leq \frac{l}{n^3} + \frac{1}{d^2}$. Hence

$\Pr[\exists l \in [1, z].s.t. W_l - C_l > (2 \sum_{a=1}^l d_a) \sqrt{2 \ln d}] \leq \frac{z^2}{n^3} + \frac{z}{d^2} \leq \frac{1}{n} + \frac{1}{d}$.

□

Now we define A_l^{OPT} and A_l in a similar manner as in the previous section. We only consider the iterations that correspond to computing $Z_{a,min}$, $a = 1, 2, \dots, z$ to define the final random deletion process.

Definition 7. A_l is defined as all the indices of the symbols that are matched or deleted by **Random-deletion** process up to and including time when processing of X_l finishes.

Definition 8. For every time $l \in [1, z]$, $A_l^{OPT} = \{i | i \in A_l \text{ or } i \text{ is matched by the optimal algorithm with some symbol with index in } A_l\}$.

We have the following corollary.

Corollary 18. For all $l \in [1, z]$, $|A_l^{OPT} \setminus A_l| \leq \sum_{a=1}^l d_a + (2 \sum_{a=1}^l d_a) \sqrt{2 \ln d}$ with probability at least $(1 - \frac{1}{n} - \frac{1}{d})$.

Proof. Proof follows from Lemma 9 and Lemma 17. \square

Lemma 19. For all $a \in [1, z]$, $StrEdit(Z_{a,min}, X_a) \leq d_a + |A_{a-1}^{OPT} \setminus A_{a-1}| + |A_a^{OPT} \setminus A_a|$.

Proof. Let $D(X_a)$ denote the symbols from X_a for which the matching open parentheses have already been deleted before processing on X_a started. Let $D'(X_a)$ denote the symbols from X_a for which the matching open parentheses are not included in $Z_{a,min}$. Let $E(Z_{a,min})$ denote open parentheses in $Z_{a,min}$ such that their matching close parentheses are in X'_a , $a' < a$, that is they are already deleted. Let $E'(Z_{a,min})$ denote open parentheses in $Z_{a,min}$ such that their matching close parentheses are in X'_a , $a' > a$, that is they are extra symbols from higher blocks.

$StrEdit(Z_{a,min}, X_a) \leq StrEdit(Z_a, X_a) + |D(X_a)| + |D'(X_a)| + |E(Z_{a,min})| + |E'(Z_{a,min})|$. Now all the indices of $D(X_a) \cup E(Z_{a,min})$ are in A_{a-1}^{OPT} , but none of them are in A_{a-1} . Hence $|D(X_a)| + |E(Z_{a,min})| \leq |A_{a-1}^{OPT} \setminus A_{a-1}|$. Also, the indices corresponding to matching close parentheses of $E'(Z_{a,min})$ and $D'(X_a)$ are in A_a^{OPT} but not in A_a . Hence $|D'(X_a)| + |E'(Z_{a,min})| \leq |A_a^{OPT} \setminus A_a|$. Therefore, the lemma follows. \square

Since the edit distance computed by the refined algorithm is at most $\sum_{a=1}^z StrEdit(Z_{a,min}, X_a)$, we get Theorem 2 using Lemma 19 and Corollary 18.

IV. FURTHER REFINEMENT: MAIN ALGORITHM & ANALYSIS

As before, we first run the process of **Random-deletion**. For each run of *Random-deletion*, the algorithm proceeds in phases with at most $\lceil \log_2 z \rceil + 1$ phases. We repeat this entire procedure $3 \log_b n$ times, $b = 1.24$ (as before) and return the minimum edit distance computed over these runs and the corresponding well-formed string. We now describe the algorithm corresponding to a single run of *random-deletion*.

Let us use $X_a^1 = X_a, Y_a^1 = Y_a$ to denote the blocks in the first phase. Consider the part of *Random-deletion* from the start of processing X_a^1 to finish either X_a^1 or Y_a^1 whichever happens first. Since this part of the random deletion (from the start of X_a^1 to the completion of either X_a^1 or Y_a^1) is contained entirely within block $Y_a^1 X_a^1$, we call this part *local¹ to block a*. Let $A_a^{local^1}$ denote the indices of all the symbols that are matched or deleted during the *local¹* steps in block a . Let $A_a^{OPT,local^1}$ be the union of $A_a^{local^1}$ and the indices of symbols that are matched with some symbol with indices in $A_a^{local^1}$ in the optimal solution. We call $A_a^{OPT,local^1} \setminus A_a^{local^1}$ the *local¹ error*, denoted $local-error^1(Y_a^1, X_a^1)$.

Create substrings L_a^1 corresponding to *local¹* moves in block a , $a = 1, \dots, z$. Compute STREDIT between $L_a^1 \cap Y_a^1$ to $L_a^1 \cap X_a^1$. Remove all these substrings from further consideration. The phase 1 ends here.

We can now view the remaining string as $Y_1^2 X_1^2 Y_2^2 X_2^2 \dots Y_z^2 X_z^2$, after deleting any prefix of

open parentheses and any suffix of close parentheses. Consider any Y_a^2, X_a^2 . Let them span the original blocks $Y_{a_1} X_{a_1} Y_{a_1+1} X_{a_1+1} \dots Y_{a_2} X_{a_2}$. Consider the part of *Random-deletion* from the start of processing X_{a_1} to the completion of either Y_{a_1} or X_{a_2} whichever happens first. Since this part of the random deletion remains confined within block $Y_a^2 X_a^2$, we call this part *local² to block a*. Let $A_a^{local^2}$ denote the indices of all the symbols that are matched or deleted during the *local²* steps in block a . Let $A_a^{OPT,local^2}$ be the union of $A_a^{local^2}$ and the indices of symbols that are matched with some symbol with indices in $A_a^{local^2}$ in the optimal solution. We call $A_a^{OPT,local^2} \setminus A_a^{local^2}$ the *local² error*, denoted $local-error^2(Y_a^2, X_a^2)$.

Create substrings L_a^2 corresponding to *local²* moves in block a , $a = 1, \dots, z^2$. Compute STREDIT between $L_a^2 \cap Y_a^2$ to $L_a^2 \cap X_a^2$. Remove all these substrings from further consideration.

We continue in this fashion until the remaining string becomes empty. We can define *localⁱ* moves, $A_a^{local^i}$, $A_a^{OPT,local^i}$ and *local-errorⁱ* (Y_a^i, X_a^i) accordingly.

Definition 9. For i th phase blocks $Y_a^i X_a^i$, if they span original blocks $Y_{a_1} X_{a_1} Y_{a_1+1} X_{a_1+1} \dots Y_{a_2} X_{a_2}$, then part of random deletion from the start of processing X_{a_1} to finish either Y_{a_1} or X_{a_2} whichever happens first, remains confined in block $Y_a^i X_a^i$ and is defined as *localⁱ move*.

Definition 10. For any $i \in \mathbb{N}$, $A_a^{local^i}$ denote the indices of all the symbols that are matched or deleted during the *localⁱ* steps in block a .

Definition 11. For any $i \in \mathbb{N}$, $A_a^{OPT,local^i}$ denote the union of $A_a^{local^i}$ and the indices of symbols matched with some symbol in $A_a^{local^i}$.

Definition 12. For any $i \in \mathbb{N}$, $A_a^{OPT,local^i} \setminus A_a^{local^i}$ is defined as the *localⁱ error*, *local-errorⁱ*.

We now summarize the algorithm below.

Algorithm:

Given the input $\sigma = Y_1 X_1 Y_2 X_2 \dots Y_z X_z$, the algorithm is as follows

- $MinEdit = \infty$,
- For $iteration = 1, iteration \leq 3 \log_b n, iteration + +$
 - Run **Random-deletion** process.
 - Set $i = 1, z^1 = z, edit = 0$, and for $a = 1, 2, \dots, z^1$, $X_a^1 = X_a, Y_a^1 = Y_a$.
 - While σ is not empty
 - * Consider the part of random-deletion from the start of processing X_a^i to finish either X_a^i or Y_a^i whichever happens first.
 - * Create substrings $L_a^i, a = 1, 2, \dots, z^i$ which correspond to *localⁱ* moves. Compute $StrEdit(L_a^i \cap Y_a^i, L_a^i \cap X_a^i)$ to match $L_a^i \cap Y_a^i$ to $L_a^i \cap X_a^i$ and add the required number of edits to $edit$.
 - * Remove $L_a^i, a = 1, 2, \dots, z^i$ from σ , write the remaining string as $Y_1^{i+1} X_1^{i+1} Y_2^{i+1} X_2^{i+1} \dots Y_{z^{i+1}}^{i+1} X_{z^{i+1}}^{i+1}$, possibly

by deleting any prefix of close parentheses and any suffix of open parentheses. The number of such deletions is also added to $edit$. Set $i = i + 1$

- End While
- If ($edit < MinEdit$) set $MinEdit = edit$
- End For

- Return $MinEdit$ as the computed edit distance.

Of course, the algorithm can compute the well-formed string by editing the parentheses that have been modified in the process through STREDIT operations.

Lemma 20. *There exists at least one iteration among $3 \log_b n$, such that for all $a' \leq b'$, (P1) the number of deletions made by random deletion process between the start of processing $X_{a'}$ and finishing either $X_{b'}$ or $Y_{a'}$ whichever happens first is at most $2d(a', b')^2$, where $d(a', b')$ is the number of deletions the optimal algorithm performs starting from the beginning of $X_{a'}$ to complete either $X_{b'}$ or $Y_{a'}$ whichever happens first with probability at least $(1 - \frac{1}{n})$.*

A. Analysis of Approximation Factor

Lemma 21. *Considering phase l which has z^l blocks, the total edit cost paid in phase l of the returned solution is $\sum_{a=1}^{z^l} StrEdit(L_a^l \cap Y_a^l, L_a^l \cap X_a^l) \leq \beta(n)(d + l * d\sqrt{24 \ln d})$ with probability at least $1 - \frac{1}{n} - \frac{1}{d}$.*

Proof. Consider the iteration in which Lemma 20 holds, that is we have property (P1). We again fix an optimal stack based algorithm and refer to it as the optimal algorithm.

Phase 1. Consider Y_a^1, X_a^1 . We know $Y_a^1 = Y_a$ and $X_a^1 = X_a$. If no symbols of either of Y_a or X_a is matched by the optimal algorithm outside of block a (that is they are matched to each other), then let d_a^1 denote the optimal edit distance to match Y_a with X_a . Consider $d_a^1 \geq 2$.

By Lemma 20, the random deletion process takes at most $2(d_a^1)^2$ steps within $local_a^1$ moves. Now from Lemma 9, local error, $A_a^{OPT, local^1} \setminus A_a^{local^1} \leq d_a^1 + W_a^{local^1} - C_a^{local^1}$ where $W_a^{local^1}$ is the number of wrong moves taken during $local^1$ steps in block a and similarly $C_a^{local^1}$ is the number of correct steps taken during the $local^1$ steps in block a . Since the total number of local steps is at most $2(d_a^1)^2$ and wrong steps are taken with probability at most $\frac{1}{2}$, hence by a standard application of the Chernoff bound or by Azuma's inequality for simple coin flips as in Lemma 17, local error, $local-error^1(Y_a^1, X_a^1) \leq d_a^1 + d_a^1 \sqrt{12 \ln d_a^1} \leq d_a^1 \sqrt{24 \ln d_a^1}$ with probability at least $1 - \frac{1}{d^3}$.

Hence $StrEdit(L_a^1 \cap Y_a^1, L_a^1 \cap X_a^1) \leq d_a + local-error^1(Y_a^1, X_a^1) \leq d_a^1 + d_a^1 \sqrt{24 \ln d_a^1}$. If $d_a^1 = 1$, then $local-error^1(Y_a^1, X_a^1) \leq 2$ and $StrEdit(L_a^1 \cap Y_a^1, L_a^1 \cap X_a^1) \leq 3d_a$.

If some suffix of X_a^1 is matched outside of block a , then let $X_a^{1,p}$ be the prefix of X_a^1 which is matched inside. Consider $Y_a^1 X_a^{1,p}$. Let d_a^1 denote the optimal edit distance to match Y_a with X_a^p . Follow the same argument as above.

Hence again $StrEdit(L_a^1 \cap Y_a^1, L_a^1 \cap X_a^1) \leq d_a + local-error^1(Y_a^1, X_a^1) \leq d_a^1 + d_a^1 \sqrt{24 \ln d_a^1}$.

Similarly if some prefix of Y_a^1 is matched outside of block a , then let $Y_a^{1,s}$ be the suffix of Y_a^1 which is matched inside, and carry out the argument considering $Y_a^{1,s} X_a^1$.

Hence due to phase 1, the total edit cost paid is at most $\sum_{a=1}^{z^1} StrEdit(L_a^1 \cap Y_a^1, L_a^1 \cap X_a^1) \leq \sum_{a=1}^{z^1} d_a^1 + d_a^1 \sqrt{24 \ln d_a^1} \leq d + d\sqrt{24 \ln d}$. This also holds when d_a^1 or d is small.

Phase 1. Let Y_a^l and X_a^l contain blocks from index $[g^{l-1}, h^{l-1}]$ of level $l-1$, all of which together contain blocks $[g^{l-2}, h^{l-2}]$ from level $l-2$ and so on, finally $[g^1, h^1] = [g, h]$ of blocks from level 1. Let d_g^h denote the optimal edit distance to match Y_a^l and X_a^l excluding the symbols that are matched outside of blocks $[g, h]$ by the optimal algorithm, and again assume $d_g^h \geq 2$. $d_g^h = 1$ is easy and can be handled exactly as in phase 1 when d_a^1 was 1.

Suppose, *Random-deletion* selects $R \subset Y_a^l$ and $T \subset X_a^l$ to match. Note that either $R = Y_a^l$ or $T = X_a^l$. Let $D(R, T)$ denote all the symbols in R, T such that their matching parentheses belong to blocks either outside of index $[g, h]$ or they exist at phase l but are not included. Let $E(R, T)$ denote all the symbols in R, T such that their matching parentheses belong to blocks $[g, h]$ but have already been deleted in phases $1, 2, \dots, l-1$.

Then $StrEdit(L_a^l \cap Y_a^l, L_a^l \cap X_a^l) \leq d_g^h + |D(R, T)| + |E(R, T)|$. Now, $|D(R, T)| \leq \frac{local-error^l(Y_a^l, X_a^l)}{A_a^{OPT, local^l} \setminus A_a^{local^l}} \leq d_g^h + d_g^h \sqrt{12 \ln d_g^h} \leq d_g^h \sqrt{20 \ln d_g^h}$.

For each $x \in E(R, T)$, consider the largest phase $\eta \in [1, 2, \dots, l-1]$ such that its matching parenthesis y existed before start of the η th phase but does not exist after the end of the η th phase. It is possible, either y belongs to the same block in phase η , say $r \in [g^\eta, h^\eta]$, or in different blocks, say r and $s \in [g^\eta, h^\eta]$. In the first case, x can be charged to y which contributes 1 to $local-error^\eta(Y_r^\eta, X_r^\eta)$. In the second case, x can again be charged to y which contributes 1 to $local-error^\eta(Y_s^\eta, X_s^\eta)$.

Hence, $|E(R, T)| \leq \sum_{j=l-1}^1 \sum_{i=g^j}^{h^j} local-error^j(Y_i^j, X_i^j)$. $StrEdit(L_a^l \cap Y_a^l, L_a^l \cap X_a^l) \leq d_g^h + local-error^l(Y_a^l, X_a^l) + \sum_{j=l-1}^1 \sum_{i=g^j}^{h^j} local-error^j(Y_i^j, X_i^j) \leq d_g^h + d_g^h \sqrt{24 \ln d_g^h} + (l-1) * d_g^h \sqrt{24 \ln d_g^h} \leq d_g^h + l d_g^h \sqrt{24 \ln d_g^h}$.

Hence due to phase l , the total edit cost paid is at most $\sum_{a=1}^{z^l} StrEdit(L_a^l \cap Y_a^l, L_a^l \cap X_a^l) \leq d + l * d\sqrt{24 \ln d}$.

Now, since we are using a $\beta(n)$ -approximation algorithm for STREDIT, we get the total edit cost paid during phase l is at most $\beta(n)(d + l * d\sqrt{24 \ln d})$. For this bound to be correct, all the $local-error$ estimates have to be correct. The number of blocks reduces by $\frac{1}{2}$ from one phase to the next. Hence, the total number of local error estimates is $\Theta(z)$. We have considered the iteration such that property (P1) holds for all sequence of blocks (see Lemma 20). Given (P1) holds, since there are a total of $\Theta(z)$ blocks over all the phases, with probability at least $1 - \frac{\Theta(z)}{d^3} > 1 - \frac{1}{d}$, all the $local-error$ bounds used in the analysis are correct. Since, (P1) holds with probability at least $(1 - \frac{1}{n})$, with probability at least $(1 - \frac{1}{n})(1 - \frac{1}{d}) > 1 - \frac{1}{n} - \frac{1}{d}$, we get the total edit cost paid

during phase l is at most $\beta(n)(d + l * d\sqrt{24 \ln d})$. \square

Noting that the number of phases can be at most $\lceil \log z \rceil + 1$, and summing the edit cost over phases, we get

Lemma 22. *The total edit cost paid is at most $O((\log z)^2 \beta(n) \sqrt{\ln d})$ with probability at least $1 - \frac{1}{n} - \frac{1}{d}$.*

Since, for a particular iteration, each STREDIT is run on a disjoint subsequence, the following theorem is obtained.

Theorem 23. *There exists an algorithm that obtains an $O(\beta(n) \sqrt{\ln d} (\log z)^2)$ -approximation factor for edit distance computation to DYCK(s) for any $s \geq 2$ in $O(n \log n + \alpha(n))$ time with probability at least $(1 - \frac{1}{n} - \frac{1}{d})$, where there exists an algorithm for STREDIT running in $\alpha(n)$ time that achieves an approximation factor of $\beta(n)$.*

A more careful charging argument improves the approximation factor to $O(\beta(n) \sqrt{\ln d} \log z)$ (Theorem 3).

Finally, the approximation factor can be further improved to $O(\beta(n) \log z)$, if we consider $O(n^\epsilon \log n)$ iterations instead of $O(\log n)$. This enables us to use the stronger Corollary 14 over Lemma 13 leading to the improvement.

Note. Due to local computations, it is possible to parallelize this algorithm.

V. FUTURE DIRECTIONS

This work raises several open questions.

(1) Is it possible to characterize the general class of grammars for which **Random-deletion** and/or its subsequent refinements can be applied?

(2) The lower bound result of Lee [24] precludes an exact algorithm (as well as one with nontrivial multiplicative approximation factor) for language edit distance problem for general context free grammars in time less than boolean matrix multiplication. Does this lower bound also hold for DYCK(s)? Is boolean matrix multiplication time enough to compute language edit distance for any arbitrary context free grammar? Is it possible to achieve nontrivial approximation guarantees when the time required for parsing a grammar is allowed?

(3) Currently there is a gap of $\log z$ in the approximation factors of string and DYCK(s) edit distance problems. Is it possible to get rid off this gap, or establish the necessity of it? What happens in other computation models such as streaming?

(4) Finally, our multi-phase random walk technique can be very useful to provide a systematic way to speed up dynamic programming algorithms for sequence alignment type problems. It will be interesting to understand whether this method can lead to a better bound for string edit distance computation.

VI. ACKNOWLEDGMENT

The author would like to thank Divesh Srivastava and Flip Korn for asking this nice question which led to their joint work [20], Arya Mazumdar for many helpful discussions and Mikkel Thorup for comments at an early stage of this work.

REFERENCES

- [1] Alfred V. Aho and Thomas G. Peterson. A minimum distance error-correcting parser for context-free languages. *SIAM J. Comput.*, 1(4), 1972.
- [2] Miklós Ajtai. The invasiveness of off-line memory checking. STOC, 2002.
- [3] Noga Alon, Michael Krivelevich, Ilan Newman, and Mario Szegedy. Regular languages are testable with a constant number of queries. *SIAM J. Comput.*, 30(6), December 2001.
- [4] Alexandr Andoni, Robert Krauthgamer, and Krzysztof Onak. Poly-logarithmic approximation for edit distance and the asymmetric query complexity. In *FOCS*. IEEE, 2010.
- [5] Alexandr Andoni and Krzysztof Onak. Approximating edit distance in near-linear time. STOC, pages 199–204, 2009.
- [6] Ziv Bar-Yossef, T. S. Jayram, Robert Krauthgamer, and Ravi Kumar. Approximating edit distance efficiently. *FOCS*, 2004.
- [7] Tuğkan Batu, Funda Ergun, and Cenk Sahinalp. Oblivious string embeddings and edit distance approximations. *SODA*, 2006.
- [8] Manuel Blum, Will Evans, Peter Gemmel, Sampath Kannan, and Moni Naor. Checking the correctness of memories. *FOCS*, 1991.
- [9] Amit Chakrabarti, Graham Cormode, Ranganath Kondapally, and Andrew McGregor. Information cost tradeoffs for augmented index and streaming language recognition. *FOCS*, 2010.
- [10] Noam Chomsky and Marcel Paul Schützenberger. *The Algebraic Theory of Context-Free Languages*. Studies in Logic. North-Holland Publishing, Amsterdam, 1963.
- [11] Matthew Chu, Sampath Kannan, and Andrew McGregor. Checking and spot-checking the correctness of priority queues. *ICALP*, 2007.
- [12] Richard Cole and Ramesh Hariharan. Approximate string matching: A simpler faster algorithm. *SIAM J. Comput.*, 31(6), June 2002.
- [13] Cynthia Dwork, Moni Naor, Guy N. Rothblum, and Vinod Vaikuntanathan. How efficient can memory checking be? *TCC*, 2009.
- [14] C. N. Fischer and J. Mauney. On the role of error productions in syntactic error correction. *Comput. Lang.*, 5(3-4), January 1980.
- [15] Charles N. Fischer and Jon Mauney. A simple, fast, and effective ll(1) error repair algorithm. *Acta Inf.*, 29(2), April 1992.
- [16] Steven Grijzenhout and Maarten Marx. The quality of the xml web. *Web Semant.*, 19, March 2013.
- [17] R.R. Gutell, J.J. Cannone, Z Shang, Y Du, and M.J. Serra. A story: unpaired adenosine bases in ribosomal rnas. *J Mol Biol*, 2010.
- [18] M. A. Harrison. *Introduction to Formal Language Theory*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1st edition, 1978.
- [19] Ik-Soon Kim and Kwang-Moo Choe. Error repair with validation in lr-based parsing. *ACM Trans. Program. Lang. Syst.*, 23(4), July 2001.
- [20] Flip Korn, Barna Saha, Divesh Srivastava, and Shanshan Ying. On repairing structural problems in semi-structured data. *VLDB*, 2013.
- [21] Dexter C. Kozen. *Automata and Computability*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 1st edition, 1997.
- [22] Andreas Krebs, Nutan Limaye, and Srikanth Srinivasan. Streaming algorithms for recognizing nearly well-parenthesized expressions. *MFCs*, 2011.
- [23] Gad M. Landau, Eugene W. Myers, and Jeanette P. Schmidt. Incremental string comparison. *SIAM J. Comput.*, 27(2), April 1998.
- [24] Lillian Lee. Fast context-free grammar parsing requires fast boolean matrix multiplication. *J. ACM*, 49(1), January 2002.
- [25] Frédéric Magniez, Claire Mathieu, and Ashwin Nayak. Recognizing well-parenthesized expressions in the streaming model. *STOC '10*, 2010.
- [26] Bruce McKenzie, Corey Yeatman, and Lorraine De Vere. Error repair in shift-reduce parsers. *ACM Transactions on Programming Languages and Systems*, 17, 1995.
- [27] Gene Myers. Approximately matching context-free languages. *Information Processing Letters*, 54, 1995.
- [28] Moni Naor and Guy N. Rothblum. The complexity of online memory checking. *J. ACM*, 56(1), February 2009.
- [29] Rafail Ostrovsky and Yuval Rabani. Low distortion embeddings for edit distance. *J. ACM*, 54(5), October 2007.
- [30] Michal Parnas, Dana Ron, and Ronitt Rubinfeld. Testing membership in parenthesis languages. *Random Struct. Algorithms*, 22(1), January 2003.
- [31] S. C. Sahinalp and U. Vishkin. Efficient approximate and dynamic matching of patterns using a labeling paradigm. *FOCS*, 1996.
- [32] Balaji Venkatachalam, Dan Gusfield, and Yelena Frid. Faster algorithms for rna-folding using the four-russians method. In *WABI*, 2013.