

D^3 as a 2-MCFL

Abstract. We discuss the open problem of parsing the Dyck language of 3 symbols, D^3 , using a 2-Multiple Context-Free Grammar. We tackle this problem by implementing a number of novel techniques and present the associated software packages we developed.

Keywords: Dyck Language; multiple context free grammars (MCFG)

1 Introduction

Our goal with this paper is the analysis of the 3-dimensional Dyck language, D^3 , under the scope of a 2-multiple context-free language, 2-MCFL. For brevity's sake, this chapter only serves as a brief introductory guide towards relevant papers, where the interested reader will find definitions, properties and various correspondences of the problem.

D^3 is defined over a lexicographically ordered alphabet (a, b, c) as the generalized language of well-balanced parentheses[4]. It is a subset of MIX[2], which has already been proven expressible by a 2-MCFG[6]; the class of multiple context-free grammars that operate on pairs of strings[1].

There are a number of interesting correspondences to D^3 . Firstly, a word of D^3 can be presented as a *standard Young Tableau*, which is a rectangular table with strictly ascending rows/columns containing as entries the numbers $\{1, 2, \dots, n\}$ where n the total number of symbols in the word. The Young Tableau can be obtained by placing (in order) each character's index to the row corresponding to its lexicographical ordering (e.g. an a is placed on the first row)[4].

Another correspondence exists between D^3 and combinatorial *Spider Webs*, a special category of directed planar graphs embedded on a disk[5]. Spider Webs can be obtained through the application of a set of rules, known as the *Growth Algorithm*, which operates on pairs of neighbouring nodes, collapsing them into a singular intermediate node, transforming them into a new pair or eliminating them altogether. At termination, this process produces a well-formed Spider Web, which, in the context of D^3 , can be interpreted as a visual representation of parsing a word.

A bijection links Young Tableaux with Spider Webs. Specifically, the *Jeu-de-taquin* algorithm can be applied on a Young Tableau, which transforms it to another one through an act called *promotion*. Subsequent promotions will eventually result in the initial tableau. Promotion thus defines an equivalence class, which we call an *orbit*[5].

2 Modeling Techniques

We now present a number of novel techniques that we developed as an attempt to solve the problem at hand, incrementally moving towards more complex and abstract grammars. For the purpose of experimentation we have implemented these techniques, based on a software library for parsing MCFGs[3]. The resulting Python code is open-source and available online¹.

2.1 Triple Insertion

To set things off, we start with the grammar of *triple insertion*. This grammar operates on non-terminals $W(x, y)$, producing $W(x', y')$ with an additional triplet a, b, c that respects the partial orders $x < y$ and $a < b < c$. The end-word is produced through the concatenation of (x, y) .

$$\begin{array}{ll}
 S(xy) \leftarrow W(x, y). & (1) \\
 W(\epsilon, xy\mathbf{abc}) \leftarrow W(x, y). & (2) \\
 \dots & \\
 W(\mathbf{abc}xy, \epsilon) \leftarrow W(x, y). & (61) \\
 W(\epsilon, \mathbf{abc}). & (62) \\
 \dots & \\
 W(\mathbf{abc}, \epsilon). & (65)
 \end{array}
 \qquad
 \begin{array}{l}
 S(xy) \leftarrow W(x, y). \\
 \mathcal{O}_2[W \leftarrow \epsilon \mid \{a < b < c\}]. \\
 \mathcal{O}_2[W \leftarrow W \mid \{x < y, a < b < c\}].
 \end{array}$$

Fig. 2. \mathcal{G}_0 : Meta-grammar of triple insertions

Fig. 1. Grammar of triple insertions

Despite being conceptually simple, this grammar consists of a large number of rules. Its expressivity is also limited; the prominent weak point is its inability to manage the effect of *straddling*, namely the generation of words whose constituents display complex interleaving patterns. Refer to Fig. 7 for an example.

2.2 Meta-Grammars

To address the issue of rule size, we introduce the notion of *meta-grammars*, loosely inspired by Van Wijngaarden's work[7], which allows a more abstract view of the grammar as a whole. Specifically, we define \mathcal{O} as the *meta-rule* which, given a rule format, a set of partial orders (over the tuple indices of its premises and/or newly added terminal symbols), and the MCFG dimensionality, automatically generates all the order-respecting permutations. An example

¹ <https://github.com/omelkonian/dyck>

of how we can abstract away from explicitly enumerating the entirety of our initial rules is showcased in Fig. 2.

This approach enhances the potential expressivity of our grammars as well. For instance, we can now extend the previous with a single meta-rule that allows two non-terminals $W(x, y)$, $W(z, w)$ to interleave with one another, producing rearranged tuple concatenations and allowing some degree of straddling to be generated:

$$\mathcal{G}_1 : \mathcal{G}_0 + \mathcal{O}_2[\mathbb{W} \leftarrow \mathbb{W}, \mathbb{W} \mid \{x < y, z < w\}].$$

The addition of this rule gets us closer to completeness, but we are still not quite there. We have thus far only used a single non-terminal, not utilizing the expressivity that an MCFG allows. To that end, we propose non-terminals to represent incomplete word *states*; that is, words that either have an extra symbol or miss one. The former are *positive* states, whereas the latter are *negative*. The inclusion of these extra states would allow for more intricate interactions. Interestingly, there is a direct correspondence between these non-terminals and the nodes of Petersen's growth algorithm.

This meta-grammar, given below, consists of base cases for positive states, possible state interactions, closures of pairs of inverse polarity and a universally quantified meta-rule that allows the combination of any incomplete state with a well-formed one (i.e. non-terminal W).

$$\begin{array}{ll} \mathbb{S}(xy) \leftarrow \mathbb{W}(x, y). & \mathcal{O}_2[\mathbb{A}^- \leftarrow \mathbb{B}^+, \mathbb{C}^+ \mid \{x < y < z < w\}]. \\ \mathcal{O}_2[\mathbb{W} \leftarrow \epsilon \mid \{a < b < c\}]. & \mathcal{O}_2[\mathbb{A}^+ \leftarrow \mathbb{C}^-, \mathbb{B}^- \mid \{x < y < z < w\}]. \\ \mathcal{O}_2[\mathbb{A}^+ \leftarrow \epsilon \mid \{a\}]. & \mathcal{O}_2[\mathbb{B}^+ \leftarrow \mathbb{C}^-, \mathbb{A}^- \mid \{x < y < z < w\}]. \\ \mathcal{O}_2[\mathbb{B}^+ \leftarrow \epsilon \mid \{b\}]. & \mathcal{O}_2[\mathbb{C}^+ \leftarrow \mathbb{B}^-, \mathbb{A}^- \mid \{x < y < z < w\}]. \\ \mathcal{O}_2[\mathbb{C}^+ \leftarrow \epsilon \mid \{c\}]. & \mathcal{O}_2[\mathbb{W} \leftarrow \mathbb{A}^+, \mathbb{A}^- \mid \{x < y < z < w\}]. \\ \mathcal{O}_2[\mathbb{C}^- \leftarrow \mathbb{A}^+, \mathbb{B}^+ \mid \{x < y < z < w\}]. & \mathcal{O}_2[\mathbb{W} \leftarrow \mathbb{C}^-, \mathbb{C}^+ \mid \{x < y < z < w\}]. \\ \mathcal{O}_2[\mathbb{B}^- \leftarrow \mathbb{A}^+, \mathbb{C}^+ \mid \{x < y < z < w\}]. & \forall \mathbb{K} \in \{\mathbb{A}^{+/-}, \mathbb{B}^{+/-}, \mathbb{C}^{+/-}\}: \\ & \mathcal{O}_2[\mathbb{K} \leftarrow \mathbb{K}, \mathbb{W} \mid \{x < y, z < w\}]. \end{array}$$

Fig. 3. \mathcal{G}_2 : Meta-grammar of incomplete states

A further extension can be achieved through universally quantifying the notion of triple insertion, which is unique in the sense that it can insert three different terminals, each at a different position:

$$\mathcal{G}_3 : \mathcal{G}_2 + \forall \mathbb{K} \in \{\mathbb{A}^{+/-}, \mathbb{B}^{+/-}, \mathbb{C}^{+/-}\} : \mathcal{O}_2[\mathbb{K} \leftarrow \mathbb{K} \mid \{x < y, a < b < c\}].$$

2.3 Rule Inference

The improved performance of the above approaches again proved insufficient to completely parse D^3 . Our meta-rules are over-constrained by imposing a total order on the tuple elements, due to their inability to keep track of where the extra character(s) is. To overcome this, we split each state into multiple position-aware, *refined* states. Doing so revealed a vast amount of new interactions, as evidenced by the below alteration to the original A^+ , B^+ interaction (where y can now occur after z or w):

$$\mathcal{O}_2[C^- \leftarrow A_{left}^+, B^+ \mid \{x < y, x < z < w\}].$$

In order to accommodate the interactions between this increased number of states, we need to keep track of both internal and external order constraints. At this point, the abstraction offered by our meta-grammar approach does not cover our needs any more. The same difficulty that we had encountered before is prominent once more, except now at an even higher level.

As a solution to the aforementioned limitation, we propose a system that can automatically create a full-blown m-MCFG given only the states it consists of. To accomplish this, we assign each state a unique *descriptor* that specifies the content of its tuple's elements. Aligning these descriptors with the tuple, we can then infer the descriptor of the resulting tuple of every possible state interaction. For the subset of those interactions whose resulting descriptor is matched with a state, we can now automatically infer the rule.

Formally, the system is initialized with a map \mathcal{D} , whose domain, $dom(\mathcal{D})$, is a set of *state identifiers* and its codomain, $codom(\mathcal{D})$, is the set of their corresponding *state descriptors* as illustrated in Fig.4.

$W \mapsto (\epsilon, \epsilon)$
 $A_l^+ \mapsto (a, \epsilon)$
 $A_r^+ \mapsto (\epsilon, a)$
 $B_l^+ \mapsto (b, \epsilon)$
 \dots
 $B_{l,r}^- \mapsto (a, c)$
 $C_l^- \mapsto (ab, \epsilon)$
 $C_r^- \mapsto (\epsilon, ab)$
 $C_{l,r}^- \mapsto (a, b)$

Fig. 4. Map \mathcal{D}

Algorithm 1 ARIS: Automatic Rule Inference System

```

procedure aris( $\mathcal{D}$ )
  for  $X \mapsto (d_1, \dots, d_n) \in \mathcal{D}$  do
    yield  $X(d_1, \dots, d_n)$ .
  for  $X, Y \in dom(\mathcal{D})^2$  do
     $(X_{ord}, Y_{ord}) \leftarrow (x < y < \dots, z < w < \dots)$ 
    for  $(d_1, \dots, d_n) \in \mathcal{O}_2[\_ \leftarrow X, Y \mid \{X_{ord}, Y_{ord}\}]$  do
      for  $S' \in eliminate((d_1, \dots, d_n), \mathcal{D})$  do
        yield  $S'(d_1, \dots, d_n) \leftarrow X, Y$ .

procedure eliminate( $((d_1, \dots, d_n), \mathcal{D})$ )
  for  $matches \in all\_abc\_triplets(d_1, \dots, d_n)$  do
    for  $i \in 0 \dots n/3$  do
      for  $S' \in remove\_abc\_triplets(matches, i)$  do
        if  $S' \in codom(\mathcal{D})$  then
          yield  $S'$ 

```

Meta-grammars accelerated the process of creating grammars, by letting us simply describe rules instead of explicitly defining them. ARIS builds upon this notion to raise the level of abstraction even further; one needs only specify a grammar's states and its descriptors, thus eliminating the need to define rules or even meta-rules.

3 Tools & Results

3.1 Grammar Utilities

We have implemented the modelling techniques described in section 2 and distributed a Python package, called **dyck**, which provides the programmer with a *domain-specific language* close to this paper's mathematical notation. To facilitate experimentation, our package includes features such as grammar selection, time measurements, word generation and soundness/completeness checking. The following example demonstrates the definition of \mathcal{G}_1 :

```
from dyck import *
G_1 = Grammar([
    ('S <- W', {(x, y)}),
    ('W', {(a, b, c)}),
    ('W <- W', {(x, y), (a, b, c)}),
    ('W <- W, W', {(x, y), (z, w)})
])
```

3.2 Visualization

As counter-examples began to grow in size and number, we realised the necessity of a visualization tool to assist us in identifying properties they may exhibit. To that end, we distribute another Python package, called **dyckviz**, which allows the simultaneous visualization of tableau-promotion and web-rotation (grouped in their corresponding equivalence classes). An example of a web as rendered by our tool is given in Fig.5.

Young tableaux in an orbit are colour-grouped by their column indices, which sheds some light on how the *jeu-de-taquin* actually influences the structure of the corresponding Dyck words. Interesting patterns have begun to emerge, which still remain to be properly investigated.

3.3 Grammar Comparisons

We display three charts, depicting the number of rules, percentage of counter-examples and computation times of each of our grammars for D_n^3 with n ranging from 2 to 6 (where n denotes the number of *abc* triplets). Even though none of our proposed grammars is complete, we observe that as grammars get more

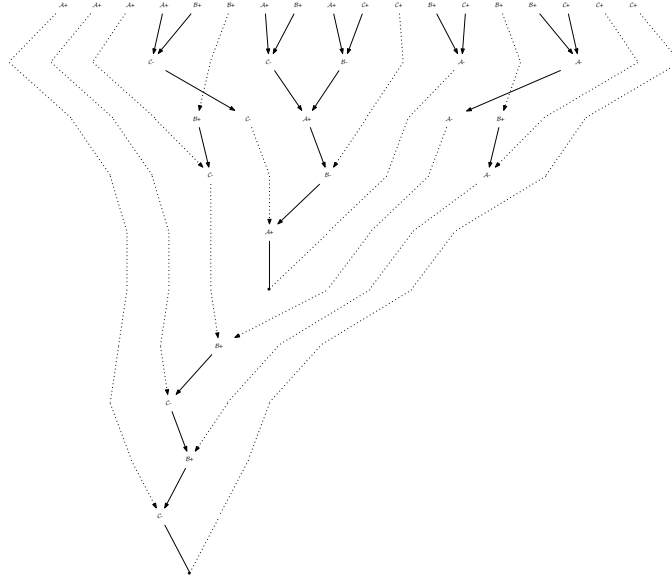


Fig. 5. Spider web of "abaacbbacbabacbcc"

abstract, the number of failing parses steadily declines. This however comes at the cost of rule size growth, which in turn is associated with an increase in computation times. What this practically means is that we are unable to continue testing more elaborate grammars or scale our results to higher orders of n (note that $\|D_n^3\|$ also has a very rapid rate of expansion)[4].

4 Road to Completeness

To our knowledge, no other attempt has come so close to modelling D^3 with a 2-MCFG. We attribute this to the combination of a pragmatic approach with results from existing theoretical work. In this chapter, we present a collection of additional ideas, which we consider worthy of further exploration.

First-Match Policy and Relinking Possibly the most intuitive way of checking whether a word w is part of D_n^3 is checking whether a pair of links occur that match a_i to b_i and b_i to $c_i \forall i \in n$. We call this process of matching the *first-match policy*. The question arises whether a grammar can accomplish inserting a triplet of a, b, c , that would abide by the first-match policy. If that were the case, it would be relatively easy to generalize this ability by induction to every $n \in \mathbb{N}$. Unfortunately, the answer is seemingly negative; the expressiveness provided by a 2-MCFG does not allow for the arbitrary insertions required. On a related

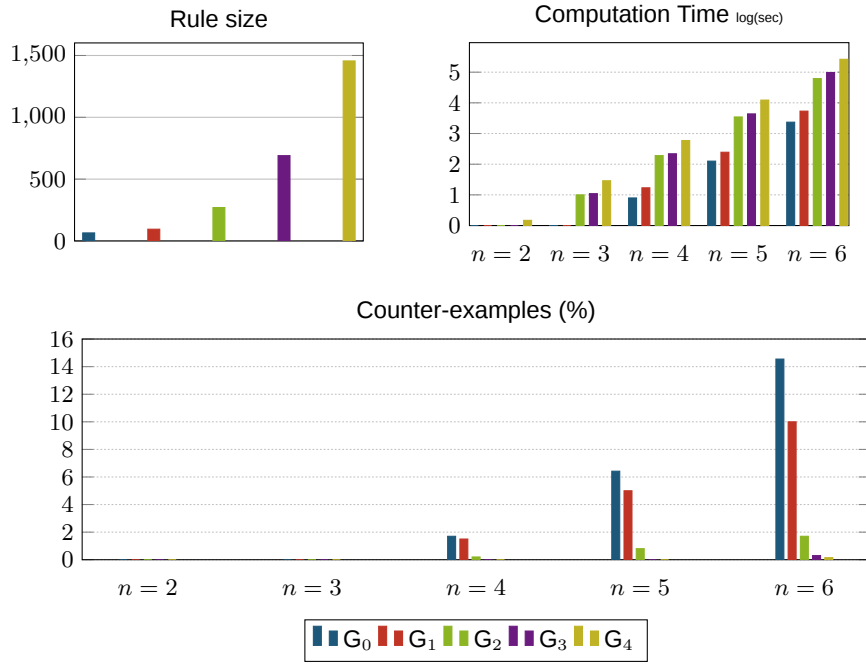


Fig. 6. Performance measures

note, being able to produce a word state $W(x, y)$ where $w = xy$ and x any possible prefix of w , gives no guarantee of being able to produce the same word with an extra triplet inserted due to the straddling property.

However, if rules existed that would allow for match-making and breaking, i.e. match *relinking*, an inserted symbol could be temporarily matched with what might be its first match-policy in a local scope, and then relink it to its correct match when merging two words together.



Fig. 7. First-match policy for "ababacbcabcc"

Growth Rules Although G_2 comes close to realizing the growth algorithm, not all of the growth rules can be translated into a 2-MCFG setting. It would be an

interesting endeavour to attempt to model the element-swapping behaviour of these rules that produce two output states, without resorting to more expressive formalisms (e.g. context-sensitive grammars).

Insights from promotion An interesting question is whether promotion can be handled by a 2-MCFG (as a *context-free rewriting system*). If so, it could be worth looking into the properties of orbits, to test for instance if there are promotions within an orbit that can be easier to solve than others. Solving a single promotion and transducing the solution to all equivalent words could then be a guideline towards completeness.

5 Conclusion

We tried to accurately present the intricacies of D^3 and the difficulties that arise when attempting to model it under the scope of a 2-MCFL. We have developed and introduced some novel techniques and tools, which we believe can be of use even outside the problem's narrow domain. We have largely expanded on the existing tools to accommodate MIX-style languages and systems of meta-grammars in general.

Despite our best efforts, the question of whether D^3 can actually be encapsulated within a 2-MCFG still remains unanswered. Regardless, this problem has been very rewarding to pursue, and we hope to have intrigued the interested reader enough to further research the subject, use our code, or strive for a solution on his/her own.

References

1. Götzmann, D.N.: Multiple context-free grammars
2. Kanazawa, M., Salvati, S.: Mix is not a tree-adjoining language. In: Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1. pp. 666–674. Association for Computational Linguistics (2012)
3. Ljunglöf, P.: Practical parsing of parallel multiple context-free grammars. In: Workshop on Tree Adjoining Grammars and Related Formalisms. p. 144 (2012)
4. Moortgat, M.: A note on multidimensional dyck languages. In: Categories and Types in Logic, Language, and Physics. pp. 279–296 (2014)
5. Petersen, T.K., Pylyavskyy, P., Rhoades, B.: Promotion and cyclic sieving via webs. *Journal of Algebraic Combinatorics* 30(1), 19–41 (2009)
6. Salvati, S.: Mix is a 2-mcfl and the word problem in z^2 is solved by a third-order collapsible pushdown automaton. *Journal of Computer and System Sciences* 81(7), 1252–1277 (2015)
7. van Wijngaarden, A.: The generative power of two-level grammars. In: ICALP (1974)