

On the Complexity and Performance of Parsing with Derivatives

Michael D. Adams
University of Utah, USA

Celeste Hollenbeck
University of Utah, USA

Matthew Might
University of Utah, USA

Abstract

Current algorithms for context-free parsing inflict a trade-off between ease of understanding, ease of implementation, theoretical complexity, and practical performance. No algorithm achieves all of these properties simultaneously.

Might et al. (2011) introduced parsing with derivatives, which handles arbitrary context-free grammars while being both easy to understand and simple to implement. Despite much initial enthusiasm and a multitude of independent implementations, its worst-case complexity has never been proven to be better than exponential. In fact, high-level arguments claiming it is fundamentally exponential have been advanced and even accepted as part of the folklore. Performance ended up being sluggish in practice, and this sluggishness was taken as informal evidence of exponentiality.

In this paper, we reexamine the performance of parsing with derivatives. We have discovered that it is not exponential but, in fact, cubic. Moreover, simple (though perhaps not obvious) modifications to the implementation by Might et al. (2011) lead to an implementation that is not only easy to understand but also highly performant in practice.

Categories and Subject Descriptors D.3.4 [Programming Languages]: Processors—Parsing

Keywords Parsing; Parsing with derivatives; Performance

1. Introduction

Although many programmers have some familiarity with parsing, few understand the intricacies of how parsing actually works. Rather than hand-write a parser, many choose to use an existing parsing tool. However, these tools are known for their maintenance and extension challenges, vague error descriptions, and frustrating shift/reduce and reduce/reduce conflicts (Merrill 1993).

In a bid to improve accessibility, Might et al. (2011) present a simple technique for parsing called parsing with derivatives (PWD). Their parser extends the Brzozowski derivative of regular expressions (Brzozowski 1964) to support context-free grammars (CFGs). It transparently handles language ambiguity and recursion and is easy to implement and understand.

PWD has been implemented in a number of languages (McGuire 2012; Vognsen 2012; Mull 2013; Shearar 2013; Byrd 2013; Engelberg 2015; Pfiel 2015). However, these tend to perform poorly, and many conjectured that the algorithm is fundamentally exponential (Cox 2010; Spiewak 2011) and could not be implemented efficiently. In fact, Might et al. (2011) report that their implementation took two seconds to parse only 31 lines of Python.

In this paper, we revisit the complexity and performance of PWD. It turns out that the run time of PWD is linearly bounded by the number of grammar nodes constructed during parsing, and we can strategically assign *unique names* to these nodes in such a way that the number of possible names is *cubic*. This means that the run time of PWD is, in fact, cubic, and the assumed exponential complexity was illusory.

Investigating further, we revisit the implementation of PWD by Might et al. (2011) by building and carefully profiling a new implementation to determine bottlenecks adversely affecting performance. We make three significant improvements over the original algorithm: *accelerated fixed points*, *improved compaction*, and *more efficient memoization*. Once these are fixed, PWD’s performance improves to match that of other general CFG parsers.

This paper makes the following contributions:

- Section 2 reviews the work by Might et al. (2011) on PWD and its key ideas.
- Section 3 investigates PWD’s complexity and shows that its upper bound is cubic instead of the exponential that was previously believed. This makes PWD’s asymptotic behavior on par with that of other general CFG parsers.
- Section 4 examines PWD’s performance and shows that targeted algorithmic improvements can achieve a speedup of almost 1000 times over the implementation in Might et al. (2011).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, or to redistribute to lists, contact the Owner/Author(s). Request permissions from permissions@acm.org or Publications Dept., ACM, Inc., fax +1 (212) 869-0481.

PLDI '16, Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation June 13 – 17, 2016, Santa Barbara, CA, USA
Copyright © 2016 held by owner/author(s). Publication rights licensed to ACM.
ACM 978-1-4503-4261-2/16/06...\$15.00
DOI: <http://dx.doi.org/10.1145/2908080.2908128>

2. Background

Brzozowski (1964) presents and Owens et al. (2009) expand upon derivatives of regular expressions as a means to recognize strings that match a given regular expression.

With PWD, Might et al. (2011) extend the concept of string recognition via Brzozowski derivatives to CFGs. The essential trick to this is handling recursively defined languages. Computing the derivative of a non-terminal may require the derivative of that same non-terminal again, causing an infinite loop. Might et al. (2011) circumvent this using a combination of memoization, laziness, and fixed points. We briefly review their technique in this section.

2.1 The Brzozowski Derivative

Brzozowski (1964) matches regular expressions against an input by successively matching each character of the input against the set of words in the semantics of that regular expression. In three steps, he computes the set of words (if any) that can validly appear after the initial input character. First, he takes the first character of the input and compares it to the first characters of the words in the semantics. Second, he keeps only the words whose first characters match and discards all others. Finally, he removes the first character from the remaining words.

Brzozowski (1964) calls this the derivative of a language and formally defines it as the following, where c is the input character and $\llbracket L \rrbracket$ is the set of words in the language L :

$$D_c(L) = \{w \mid cw \in \llbracket L \rrbracket\}$$

For example, with respect to the character f , the derivative of the language for which $\llbracket L \rrbracket = \{\text{foo}, \text{frak}, \text{bar}\}$ is $D_f(L) = \{\text{oo}, \text{rak}\}$. Because foo and frak start with the character f and bar does not, we keep only foo and frak and then remove their initial characters, leaving oo and rak .

We repeat this process with each character in the input until it is exhausted. If, after every derivative has been performed, the resulting set of words contains the empty word, ϵ , then there is some word in the original language consisting of exactly the input characters, and the language accepts the input. All of this processing takes place at parse time, so there is no parser-generation phase.

2.2 Parsing Expressions

Explicitly enumerating the possibly infinite set of words in a language can be cumbersome, so we express regular languages using the expression forms in Figure 1. For the most part, these consist of the traditional regular expression forms. The ϵ_s form is the language of the empty string, \emptyset is the empty language, c is a single token, (\circ) concatenates, and (\cup) forms alternatives. In Might et al. (2011), every expression also produces an abstract syntax tree (AST) upon success. So, ϵ_s is annotated with a subscript s indicating the AST to be returned, and the reduction form $L \hookrightarrow f$ behaves like L , except that it returns the result of applying f to the

Forms

$$\begin{aligned} L ::= & \emptyset \mid \epsilon_s \mid c \mid L_1 \circ L_2 \mid L_1 \cup L_2 \mid L \hookrightarrow f \\ s, t \in & T && \text{Abstract syntax trees} \\ f \in & T \rightarrow T && \text{Reduction functions} \end{aligned}$$

Semantics

$$\begin{aligned} \llbracket L \rrbracket & \in \wp(\Sigma^* \times T) \\ \llbracket \emptyset \rrbracket & = \{\} && \text{Empty Lang.} \\ \llbracket \epsilon_s \rrbracket & = \{(\epsilon, s)\} && \text{Empty Word} \\ \llbracket c \rrbracket & = \{(c, c)\} && \text{Token} \\ \llbracket L_1 \circ L_2 \rrbracket & = \{(uv, (s, t)) \mid (u, s) \in \llbracket L_1 \rrbracket \text{ and } (v, t) \in \llbracket L_2 \rrbracket\} && \text{Concatenation} \\ \llbracket L_1 \cup L_2 \rrbracket & = \{(u, s) \mid (u, s) \in \llbracket L_1 \rrbracket \text{ or } (u, s) \in \llbracket L_2 \rrbracket\} && \text{Alternation} \\ \llbracket L \hookrightarrow f \rrbracket & = \{(w, fs) \mid (w, s) \in \llbracket L \rrbracket\} && \text{Reduction} \end{aligned}$$

Figure 1. Parsing expression forms

AST returned by L . The semantics of these forms are as in Figure 1 and are defined as sets of accepted strings paired with the AST that returns for that string. For the purposes of parsing single tokens, c , and concatenations, (\circ) , we assume the type of ASTs includes tokens and pairs of ASTs.

Note that in this paper, we use ϵ for the empty word and ϵ_s for the parsing expression that represents a language containing only the empty word. Similarly, we use c to refer to either the single-token word or the parsing expression signifying a language containing only one token.

Also, although Might et al. (2011) include a form for Kleene star, we omit this. Once these forms are extended from regular expressions to CFGs in Section 2.5, any use of Kleene star can be replaced with a definition like the following.

$$L^* = \epsilon_s \cup (L \circ L^*)$$

2.3 Derivatives of Parsing Expressions

The derivatives of the language forms in Figure 1 with respect to a token c are shown in Figure 2. The derivative of \emptyset is \emptyset , as $\llbracket \emptyset \rrbracket$ contains no words beginning with any character. For the same reason, the derivative of ϵ_s is also \emptyset . The derivative of a token c depends on whether the input token matches c ; the result is ϵ_c if the input token matches and \emptyset if not. The derivatives of $L_1 \cup L_2$ and $L \hookrightarrow f$ merely take the derivatives of their children.

The derivative of $L_1 \circ L_2$ has two cases, depending on whether $\llbracket L_1 \rrbracket$ contains ϵ . If $\llbracket L_1 \rrbracket$ does not contain ϵ , every word in the concatenated language starts with a non-empty word from L_1 . This means the derivative of $L_1 \circ L_2$ filters and removes the first token from the words in L_1 while

$$\begin{aligned}
D_c(\emptyset) &= \emptyset \\
D_c(\epsilon) &= \emptyset \\
D_c(c') &= \begin{cases} \epsilon_c & \text{if } c = c' \\ \emptyset & \text{if } c \neq c' \end{cases} \\
D_c(L_1 \cup L_2) &= D_c(L_1) \cup D_c(L_2) \\
D_c(L_1 \circ L_2) &= \begin{cases} D_c(L_1) \circ L_2 & \text{if } \epsilon \notin \llbracket L_1 \rrbracket \\ (D_c(L_1) \circ L_2) \cup D_c(L_2) & \text{if } \epsilon \in \llbracket L_1 \rrbracket \end{cases} \\
D_c(L \hookrightarrow f) &= D_c(L) \hookrightarrow f
\end{aligned}$$

Figure 2. Derivatives of parsing expression forms

$$\begin{aligned}
\delta(\emptyset) &= \text{false} \\
\delta(\epsilon_s) &= \text{true} \\
\delta(c) &= \text{false} \\
\delta(L_1 \cup L_2) &= \delta(L_1) \text{ or } \delta(L_2) \\
\delta(L_1 \circ L_2) &= \delta(L_1) \text{ and } \delta(L_2) \\
\delta(L \hookrightarrow f) &= \delta(L)
\end{aligned}$$

Figure 3. Nullability of parsing expression forms

leaving L_2 alone. Thus, the derivative of $L_1 \circ L_2$ if L_1 does not contain ϵ is $D_c(L_1) \circ L_2$.

On the other hand, if $\llbracket L_1 \rrbracket$ *does* contain ϵ , then the derivative contains not only all the words in $D_c(L_1) \circ L_2$ but also derivatives for when the ϵ in L_1 is concatenated with words in L_2 . Since these concatenations are all words from L_2 , this adds $D_c(L_2)$ to the derivative. In this case, $D_c(L_1 \circ L_2)$ is therefore $(D_c(L_1) \circ L_2) \cup D_c(L_2)$.

2.4 Nullability

Because the derivative of a concatenation $L_1 \circ L_2$ depends on whether $\llbracket L_1 \rrbracket$ contains the empty string, ϵ , we define a nullability function, $\delta(L)$, in Figure 3 such that it returns boolean true or false when $\llbracket L \rrbracket$ respectively contains ϵ or does not. The null language, \emptyset , contains nothing, and the single-token language, c , contains only the word consisting of the token c . Because neither of these languages contain ϵ , their nullability is false. Conversely, the ϵ_s language contains only the ϵ word, so its nullability is true. The union of two languages contains ϵ if either of its children contains ϵ , so the union is nullable if either L_1 or L_2 is nullable. Given how the semantics of the concatenation $L_1 \circ L_2$ are defined in Figure 1, in order for $L_1 \circ L_2$ to contain ϵ , there must exist a uv equal to ϵ . This happens only when u and v are both ϵ , so a concatenation is nullable if and only if both its children are nullable. Finally, the words in a reduction $L \hookrightarrow f$ are those words in L , so its nullability is the nullability of L .

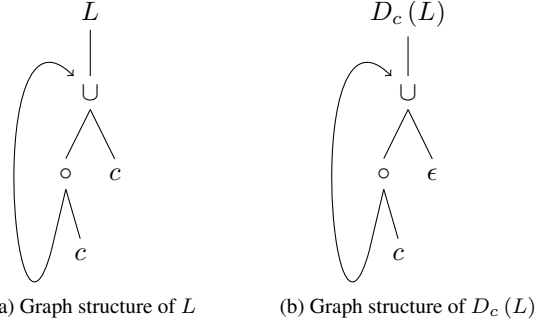


Figure 4. An example grammar and its derivative

2.5 Derivatives of Context-free Languages

2.5.1 Representation

Might et al. (2011) generalize from taking derivatives of regular expressions to taking derivatives of full CFGs. In so doing, Might et al. (2011) do not use typical CFGs but do use an equivalent construction.

First, instead of non-terminals mapping to zero or more sequences of terminals and non-terminals, they map to a parsing expression. This is akin to the parsing expressions in Ford (2004). For example, any CFG can be converted to this expression form by converting productions of the form

$$N ::= X_{11} \cdots X_{1m_1} \mid \cdots \mid X_{n1} \cdots X_{nm_n}$$

to

$$N = X_{11} \circ \dots \circ X_{1m_1} \cup \dots \cup X_{n1} \circ \dots \circ X_{nm_n}$$

Second, in the data structures representing grammars, instead of using explicitly named non-terminals, parsing expressions point directly to the non-terminal's parsing expression. For example, we may have a grammar like the following, where c is some token.

$$L = (L \circ c) \cup c$$

Might et al. (2011) represent this as the data structure in Figure 4a with the edge where L refers back to itself, forming a cycle in the data structure. For the purposes of discussion, though, we will refer to non-terminals and their names even though the actual representation uses direct pointers instead of non-terminal names.

2.5.2 Computation

A complication of this representation occurs when taking a derivative. If we blindly follow the rules in Figure 2, then the derivative of L by c is the following.

$$D_c(L) = (D_c(L) \circ c) \cup \epsilon$$

This $D_c(L)$ is recursive, so to compute $D_c(L)$, we must already know $D_c(L)$!

Might et al. (2011) solve this problem with two measures. First, they memoize their derivation function, `derive`, by keeping a table containing, for each set of arguments, the results that it returns. When `derive` is called, if the table already contains an entry for its arguments, `derive` uses the result in the entry instead of re-computing the derivative. Otherwise, `derive` performs the calculation as usual and, before returning, stores its result in the memoization table so it can be used by any further calls with those same arguments. If the same derivative is needed multiple times, this ensures it is computed only once.

On its own, memoization does not prevent infinite loops due to cycles, however, because `derive` adds memoization table entries only after it finishes computing. This is where a second measure comes into play. Before doing any recursive calls, `derive` puts a partially constructed grammar node that is missing its children into the memoization table. For the example of $D_c(L)$, we know without having to recur into L 's children that the resultant node is a \cup . Thus, we can place such a node in the memoization table before computing its children and temporarily mark its children as unknown. Any recursive calls to $D_c(L)$ can find and use this memoized result even though the derivatives of its children have not yet been calculated. When the derivatives for the node's children return, we update the children of the output node to point to the results of those derivatives. This process can be viewed as a sort of lazy computation and results in a graph structure like in Figure 4b.

Like with the derivative, computing nullability must also deal with cycles in grammars. However, memoization alone is not sufficient here. A cycle means the derivative of some node must point to one of its ancestors. With nullability, though, we must not only compute the nullability of an ancestor but also inspect its value so we can compute the nullability of the current node. This turns nullability into a least fixed point problem over the lattice of booleans. Might et al. (2011) implement this with a naive algorithm that initially assumes all nodes are not nullable and then recomputes the nullability of all nodes reachable from a particular root node, using the current values for each node. If, in the process, any nodes are newly discovered to be nullable, then all reachable nodes are re-traversed and this process is repeated until there are no more changes.

2.6 Performance

Despite PWD's simplicity and elegance, Might et al. (2011) report significant problems with its performance. Firstly, they compute a worst-case bound of $O(2^{2n}G^2)$ for a grammar of size G and an input of size n . Despite this, they note that average parse time seems to be linear in the length of the input. Unfortunately, even with this apparent linear behavior, their parser is exceedingly slow. For example, they report that a 31-line Python file took three minutes to parse! Using an optimization they call compaction that prunes branches of the derived grammars as they emerge, they report that execu-

tion time for the 31-line input comes down to two seconds. Still, this is exceedingly slow for such a small input.

3. Complexity Analysis

Might et al. (2011) report an exponential bound for their algorithm, but they never show it is a tight bound. On the contrary, it turns out that PWD can, in fact, be implemented in cubic time.

As mentioned before, at its core, PWD involves four recursive functions: `nullable?`, `derive`, `parse-null`, and `parse`. The `nullable?` and `derive` functions implement $\delta(L)$ and $D_c(L)$, respectively; the `parse-null` function extracts the final AST; and `parse` implements the outer loop over input tokens. In Section 3.1, we observe that the running times of these functions are bounded by the number of grammar nodes in the initial grammar plus the number of grammar nodes constructed during parsing. Next, in Section 3.2 we discover that the total number of nodes constructed during parsing is $O(Gn^3)$, where G is the size of the initial grammar and n is the length of the input. How to structure this part of the proof is the essential insight in our analysis and is based on counting unique names that we assign to nodes. When combined with the results from Section 3.1, this then leads to a cubic bound on the total runtime.

Throughout this section, let G be the number of grammar nodes in the initial grammar, let g be the number of nodes created during parsing, and let n be the length of the input. Also, when analyzing a memoized function, we consider the cost of the check to see if a memoized result exists for a particular input to be part of the running time of the caller instead of the callee.

3.1 Total Running Time in Terms of Grammar Nodes

First, consider `nullable?`, which computes a boolean value for each parse node in terms of a least fixed point. The implementation by Might et al. (2011) iteratively re-traverses the grammar until no new nodes can be proven `nullable?`. Such an algorithm is quadratic in the number of nodes over which `nullable?` is being computed because each traversal might update only one node. However, a more intelligent algorithm that tracks dependencies between nodes and operates over the boolean lattice can implement this function in linear time, as shown in the following lemma.

Lemma 1. *The sum of the running times of all invocations of `nullable?` is $O(G + g)$.*

Proof. The fixed point to calculate `nullable?` can be implemented by a data-flow style algorithm (Kildall 1973) that tracks which nodes need their nullability reconsidered when a given node is discovered to be nullable. Such an algorithm is linear in the product of the height of the lattice for the value stored at each node and the number of direct dependencies between nodes. In this case, the lattice is over booleans and is of constant height. Since each node directly depends

on at most two children, the number of dependencies is bounded by twice the number of nodes ever created. \square

Next, we have `derive`. Since `derive` is memoized, it is tempting to analyze it in terms of the nodes passed to it. However, each node may have its derivative taken with multiple different input tokens. The work done by `derive` thus depends on the number of tokens by which each node is derived, so we can instead simplify things by analyzing `derive` in terms of the nodes that it constructs.

Lemma 2. *The sum of the running times of all invocations of `derive` is $O(G + g)$.*

Proof. Every call to `derive` that is not cached by memoization creates at least one new node and, excluding the cost of recursive calls, does $O(1)$ work. As a result, the number of nodes created, g , is at least as great as the amount of work done. Thus, the work done by all calls to `derive` is $O(g)$ plus the work done by `nullable?`. By Lemma 1, this totals to $O(G + g)$. \square

Next, we have `parse-null`. For this part of the proof, we assume that ASTs use ambiguity nodes and a potentially cyclic graph representation. This is a common and widely used assumption when analyzing parsing algorithms. For example, algorithms like GLR (Lang 1974) and Earley (Earley 1968, 1970) are considered cubic, but only when making such assumptions. Without ambiguity nodes, the grammar $S \rightarrow S S \mid a \mid b$ has an exponential number of unique parses for strings of length n that have no repeated substrings of length greater than $\log_2 n$. Many of those parses share common sub-trees, so it does not take exponential space when represented with ambiguity nodes. Our implementation is capable of operating either with or without such a representation, but the complexity result holds only with the assumption.

Under these assumptions, `parse-null` is a simple memoized function over grammar nodes and thus is linear.

Lemma 3. *The sum of the running times of all invocations of `parse-null` is $O(G + g)$.*

Proof. Every call to `parse-null` that is not cached by memoization does $O(1)$ work, excluding the cost of recursive calls. There are at most $G + g$ such non-cached calls. \square

Finally, we have the total running time of `parse`.

Theorem 4. *The total running time of `parse` is $O(G + g)$.*

Proof. The `parse` function calls `derive` for each input token and, at the end, calls `parse-null` once. By Lemma 2 and Lemma 3, these together total $O(G + g)$. \square

3.2 Grammar Nodes in Terms of Input Length

All of the results in Section 3.1 depend on g , the number of grammar nodes created during parsing. If we look at the definition of $D_c(L)$ (i.e., `derive`) in Figure 2, most of the clauses construct only a single node and use the children of the input node only once each. When combined with memoization, for a given input token, these clauses create at most the same number of nodes as there are in the grammar for the result of the derivative just before parsing that input token. On their own, these clauses thus lead to the construction of only Gn nodes.

However, the clause for a sequence node $L_1 \circ L_2$, when L_1 is nullable, uses L_2 twice. This duplication is what led many to believe PWD was exponential; and indeed, without memoization, it would be. In order to examine this more closely, we assign unique names to each node. We choose these names such that each name is unique to the derivative of a particular node with respect to a particular token. Thus, the naming scheme matches the memoization strategy, and the memoization of `derive` ensures that two nodes with the same name are always actually the same node.

Definition 5. We give each node a unique name that is a string of symbols determined by the following rules.

Rule 5a: Nodes in the initial grammar are given a name consisting of a single unique symbol distinct from that of any other node in the grammar.

Rule 5b: When the node passed to `derive` has the name w and is a \circ node containing a nullable left child, the \cup node created by `derive` is given the name $w \bullet c$ where \bullet is a distinguished symbol that we use for this purpose and c is the token passed to `derive`.

Rule 5c: Any other node created by `derive` is given a name of the form wc , where w and c are respectively the name of the node passed to `derive` and the token passed to `derive`.

A \circ node with a nullable left child has the special case of Rule 5b because it is the only case where `derive` produces more than one node, and we need to give these nodes distinct names. These resultant nodes are a \cup node and a \circ node that is the left child of the \cup node. The introduction of the \bullet symbol in the name of the \cup node keeps this name distinct from the name of the \circ node.

As an example of these rules, Figure 5 shows the nodes and corresponding names for the nodes created when parsing the following grammar.

$$L = (L \circ L) \cup c$$

In this example, c accepts any token; the initial names are L , M , and N ; and the input is $c_1 c_2 c_3 c_4$. Each node in Figure 5 is labeled with its name in a subscript, and children that point to already existing nodes are represented with a box containing the name of that node. For example, the root of the first grammar contains the node named L as

its root, and the node named M in that tree has L as both its children. The dotted arrows in this diagram show where concatenation causes duplication. The node M produces Mc_1 , Mc_1 produces $Mc_1 \bullet c_2$ and $Mc_1 c_2$, and so on.

A nice property of these rules can be seen if we consider node names with their initial unique symbol and any \bullet symbols removed. The remaining symbols are all tokens from the input. Furthermore, these symbols are added by successive calls to `derive` and thus are substrings of the input. This lets us prove the following lemma.

Lemma 6. *The number of strings of symbols consisting of node names with their initial unique symbols and any \bullet symbols removed is $O(n^2)$.*

Proof. These strings are all substrings of the input. Flaxman et al. (2004) count the number of such substrings and show that, unsurprisingly, it is $O(n^2)$. At an intuitive level, this is because the number of positions where these substrings can start and end in the input are both linear in n . \square

In Figure 5, this can be seen by the fact that the c_1 , c_2 , c_3 , and c_4 occurring in node names are always in increasing, consecutive ranges such as $c_1 c_2 c_3$ in $Mc_1 c_2 \bullet c_3$ or $c_2 c_3$ in $Nc_2 c_3$.

Another nice property of names is that they all contain at most one occurrence of \bullet . This turns out to be critical. At an intuitive level, this implies that the \cup node involved in a duplication caused by a \circ node is never involved in another duplication.

Lemma 7. *Each node name contains at most one occurrence of the \bullet symbol.*

Proof. According to Rule 5b, a \bullet symbol is put in the name of only those \cup nodes that come from taking the derivative of a \circ node. Further derivatives of these \cup nodes can produce only more \cup nodes, so Rule 5b, which applies only to \circ nodes, cannot apply to any further derivatives of those \cup nodes. Thus, once a \bullet symbol is added to a name, another one cannot be added to the name. \square

This property can be seen in Figure 5 where no name contains more than one \bullet , and every node that does contain \bullet is a \cup node.

This then implies that every name is either of the form Nw or $Nu \bullet v$, where N is the name of an initial grammar node and both w and uv are substrings of the input. As a result, we can bound the number of possible names with the following theorem.

Theorem 8. *The total number of nodes constructed during parsing is $O(Gn^3)$.*

Proof. In a name of the form Nw or $Nu \bullet v$, the number of possible symbols for N is the size of the initial grammar, G . Also, the number of possible words for w or uv is bounded by the number of unique subwords in the input, which is

$O(n^2)$. Finally, the number of positions at which \bullet may occur within those subwords is $O(n)$. The number of unique names, and consequently the number of nodes created during parsing, is the product of these: $O(Gn^3)$. \square

3.3 Running Time in Terms of Input Length

Finally, we can conclude that the running time of parsing is cubic in the length of the input.

Theorem 9. *The running time of parse is $O(Gn^3)$.*

Proof. Use $O(Gn^3)$ in Theorem 8 for g in Theorem 4. \square

Note that this analysis does not assume the use of the process that Might et al. (2011) call compaction. Nevertheless, it does hold, in that case, if compaction rules are applied only when a node is constructed and only locally at the node being constructed. The extra cost of compaction is thus bounded by the number of nodes constructed, and compaction only ever reduces the number of nodes constructed by other parts of the parser.

4. Improving Performance in Practice

Given that PWD has a cubic running time instead of the exponential conjectured in Might et al. (2011), the question remains of why their implementation performed so poorly and whether it can be implemented more efficiently. To investigate this, we reimplemented PWD from scratch and built up the implementation one part at a time. We measured the running time as each part was added and adjusted our implementation whenever a newly added part significantly slowed down the implementation. Section 4.1 reports the final performance of the resulting parser. Aside from low-level needs to choose efficient data structures, we found three major algorithmic improvements, which are discussed in Section 4.2, Section 4.3, and Section 4.4.

The resulting parser implementation remains rather simple and easily read. The optimization of the fixed-point computation for `nullable?` (Section 4.2) takes 24 lines of Racket code, including all helpers. Compaction (Section 4.3) is implemented using smart constructors for each form that test if they are constructing a form that can be reduced. This takes 50 lines of code due to each constructor needing to have a clause for each child constructor with which it could reduce. Finally, single-entry memoization (Section 4.4) requires changing only the helpers that implement memoization, which does not increase the complexity or size of the resulting code. With all of these optimizations implemented, the core code is 62 lines of Racket code with an additional 76 lines of code for helpers.

The complete implementation can be downloaded from:

<http://www.bitbucket.com/ucombinator/derp-3>

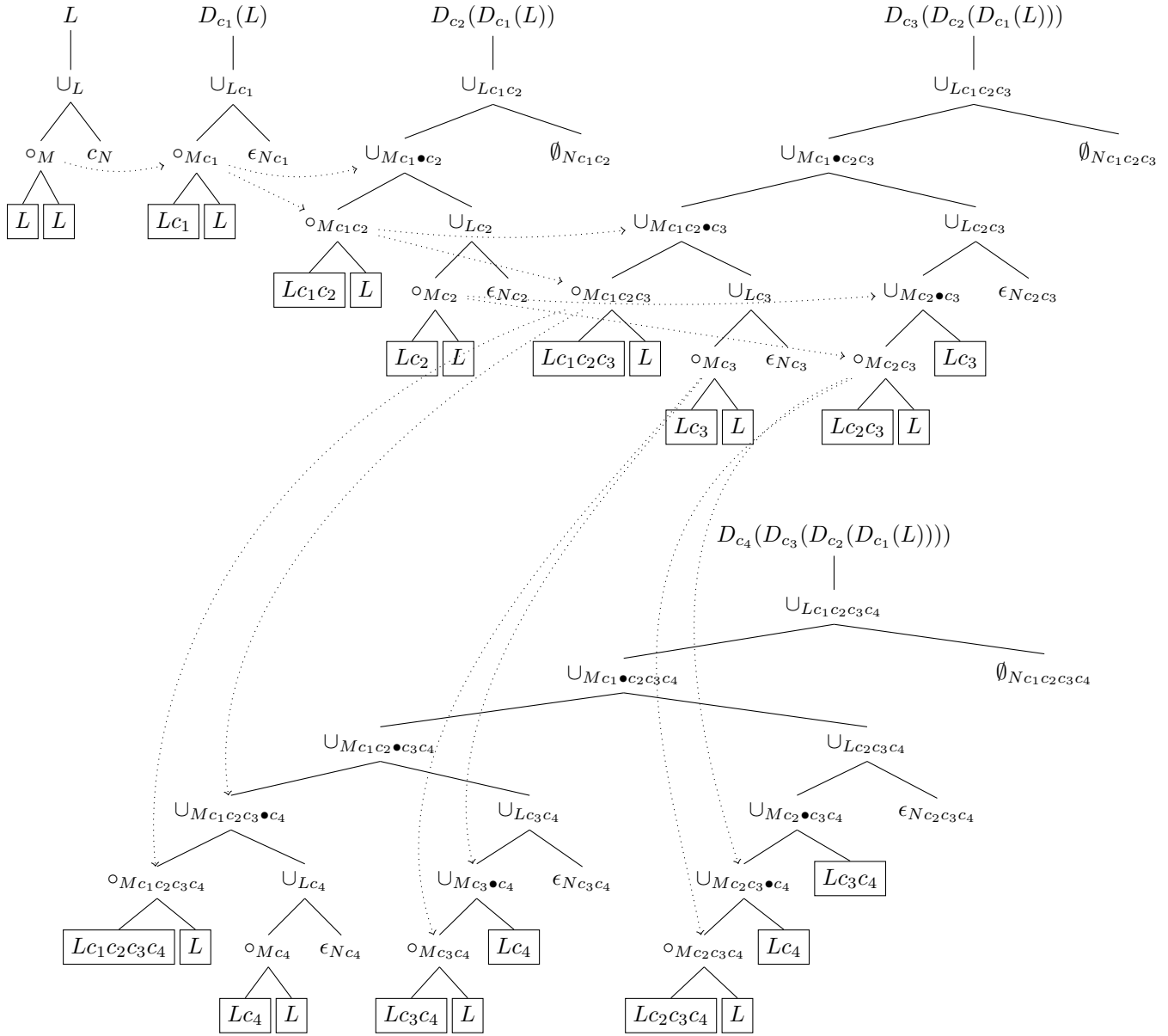


Figure 5. Worst-case behavior of PWD. Nodes are annotated with their names in subscripts.

4.1 Benchmarks

In order to test the performance of our implementation of PWD, we ran our parser on the files in the Python Standard Library version 3.4.3 (Python Software Foundation 2015a) using a grammar derived from the Python 3.4.3 specification (Python Software Foundation 2015b). The Python Standard Library includes 663 Python files, which have sizes of up to 26,125 tokens.

We compared our parser against three parsers. The first one used the original PWD implementation (Might 2013). The second one used the `parser-tools/cfg-parser` li-

brary (Parser Tools) that comes with Racket 6.1.1 (Racket). The third one used Bison version 3.0.2 (Bison).

In order to have a fair comparison against the original PWD implementation, our parser was written in Racket. For compatibility with `parser-tools/cfg-parser` and Bison, we modified our grammar to use traditional CFG productions instead of the nested parsing expressions supported by PWD and used by the Python grammar specification. The resulting grammar contained 722 productions.

The `parser-tools/cfg-parser` library uses a variant of the Earley parsing algorithm (Earley 1968, 1970), so it may not perform as well as other GLR parsers (Lang 1974).

Nevertheless, we used it because we were not able to locate a suitable GLR parser for Racket.

In order to compare against a more practical GLR parser, we included a Bison-based parser. We ran Bison in GLR mode, as the grammar resulted in 92 shift/reduce and 4 reduce/reduce conflicts. However, as the Bison-based parser is written in C and the improved PWD parser is written in Racket, the Bison-based parser has an extra performance boost that the improved PWD implementation does not have.

We ran the tests with Racket 6.1.1 and GCC 4.9.2 on a 64-bit, 2.10 GHz Intel Core i3-2310M running Ubuntu 15.04. Programs were limited to 8000 MB of RAM via `ulimit`. We tokenized files in advance and loaded those tokens into memory before benchmarking started, so only parsing time was measured when benchmarking. For each file, we computed the average of ten rounds of benchmarking that were run after at least three warm-up rounds. However, the original PWD was so slow that we could only do three rounds of benchmarking for that implementation. Each round parsed the contents of the file multiple times, so the run time lasted at least one second to avoid issues with clock quantization. We cleared memoization tables before the start of each parse. A small number of files exceeded 8000 MB of RAM when parsed by the original PWD or `parser-tools/cfg-parser` and were terminated early. We omit the results for those parsers with those files. This did not happen with the improved PWD and Bison, and the results from those parsers on those files are included. The final results are presented in Figure 6 and are normalized to measure parse time per input token.

As reported in Might et al. (2011), PWD appears to run in linear time, in practice, with a constant time per token. However, our improved parser runs on average 951 times faster than that by Might et al. (2011). It even runs 64.6 times faster than the parser that uses the `parser-tools/cfg-parser` library. As expected, our implementation ran slower than the Bison-based parser, but by only a factor of 25.2. This is quite good, considering how simple our implementation is and the differences in the implementations’ languages. We suspect that further speedups could be achieved with a more efficient implementation language.

In the remainder of this section, we explain the main high-level algorithmic techniques we discovered that achieve this performance.

4.2 Computing Fixed Points

The `nullable?` function is defined in terms of a least fixed point. The implementation in Might et al. (2011) computes this by repeatedly traversing over all grammar nodes. If the computed nullability of any node changes during that traversal, all of the nodes are traversed again. This continues until there are no more changes.

This is a fairly naive method of computing a fixed point and is quadratic in the number of nodes in the grammar as each re-traversal may update only one node that then trig-

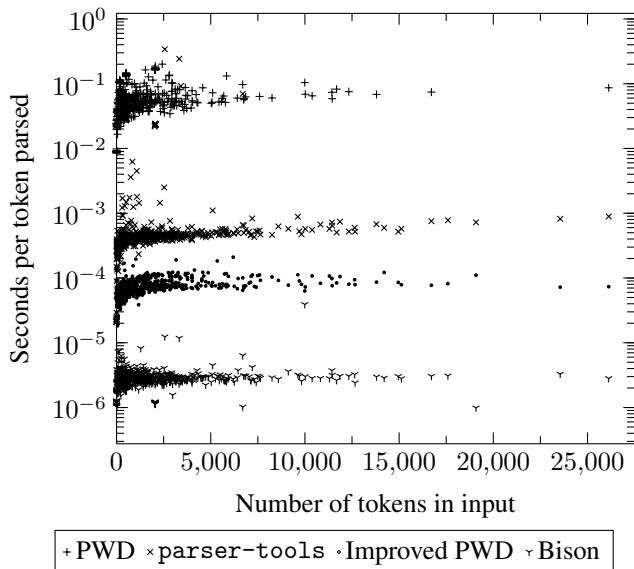


Figure 6. Performance of various parsers

gers another re-traversal. A more efficient method uses ideas from data-flow analysis (Kildall 1973) and tracks which nodes depend on which others. When the computed nullability of a node changes, only those nodes that depend on that node are revisited.

While the tracking of dependencies does incur an overhead, we can minimize this by tracking dependencies only after discovering cycles that prevent us from immediately computing the result. In other cases, we directly compute nullability with a simple recursive traversal.

We can further improve the performance of `nullable?` by distinguishing between nodes that are definitely not nullable and those that are merely assumed to be not nullable because the fixed point has not yet shown them to be nullable.

Assumed-not-nullable and definitely-not-nullable nodes behave almost exactly alike except that we may re-traverse assumed-not-nullable nodes but never re-traverse definitely-not-nullable nodes. This is because definitely-not-nullable nodes have their final value, while assumed-not-nullable nodes might not.

In many types of fixed-point problems, this is not an important distinction because there is usually no way to distinguish between these types of nodes. However, when computing nullability, we can take advantage of this because the computation of nullability is not done only once. Rather, it is called multiple times on different nodes by different executions of `derive`. Within each of these fixed points, only nodes reachable from the node passed to the initial call to `nullable?` by `derive` have their nullability computed. Later calls to `nullable?` may examine different nodes, but when they examine nodes already examined in a previous call to `nullable?` from `derive`, they can reuse information

from that previous call. Specifically, not only are nodes that are discovered by previous fixed points to be nullable still nullable, but nodes that are assumed-not-nullable at the end of a previous fixed-point calculation are now definitely-not-nullable. This is because the nodes that could cause them to be nullable are already at a value that is a fixed point and will not change due to further fixed-point calculations.

We take advantage of this by marking nodes visited by `nullable?` with a label that is unique to the call in `derive` that started the nullability computation. Then, any nodes still assumed-not-nullable that are marked with a label from a previous call are treated as definitely-not-nullable.

The end result of these optimizations is a significant reduction in the number of calls to `nullable?`. In Figure 7, we plot the number of calls to `nullable?` in our implementation relative to that of Might et al. (2011). On average, the new implementation has only 1.5% of the calls to `nullable` as that of Might et al. (2011).

4.3 Compaction

Might et al. (2011) report that a process that they call compaction improves the performance of parsing by a factor of about 90. We found similar results in our implementation, and the benchmarks in Figure 6 use compaction. However, we also discovered improvements to this process.

First, we keep the following reduction rules from Might et al. (2011) with no changes. The first three rules take advantage of the fact that \emptyset is the identity of \cup and the annihilator of \circ . The last three rules move the operations involved in producing an AST out of the way to expose the underlying grammar nodes.

$$\begin{aligned} \emptyset \cup p &\Rightarrow p \\ p \cup \emptyset &\Rightarrow p \\ \emptyset \circ p &\Rightarrow \emptyset \\ \epsilon_s \circ p &\Rightarrow p \hookrightarrow \lambda u. (s, u) \\ \epsilon_s \hookrightarrow f &\Rightarrow \epsilon_{(f s)} \\ (p \hookrightarrow f) \hookrightarrow g &\Rightarrow p \hookrightarrow (g \circ f) \end{aligned}$$

To these rules, we add the following reductions, which were overlooked in Might et al. (2011).

$$\begin{aligned} \emptyset \hookrightarrow f &\Rightarrow \emptyset \\ \epsilon_{s_1} \cup \epsilon_{s_2} &\Rightarrow \epsilon_{s_1 \cup s_2} \end{aligned}$$

We also omit the following reduction used by Might et al. (2011), as it is covered by the reductions for $\epsilon_s \circ p$ and $(p \hookrightarrow f) \hookrightarrow g$.

$$(\epsilon_s \circ p) \hookrightarrow f \Rightarrow p \hookrightarrow \lambda u. f (u, s)$$

The reader may notice that these laws are very similar to the laws for Kleene algebras (Kozen 1994). If we ignore the generated parse trees and consider only string recognition, parsing expressions are Kleene algebras. The identities for compaction have one subtle difference from those

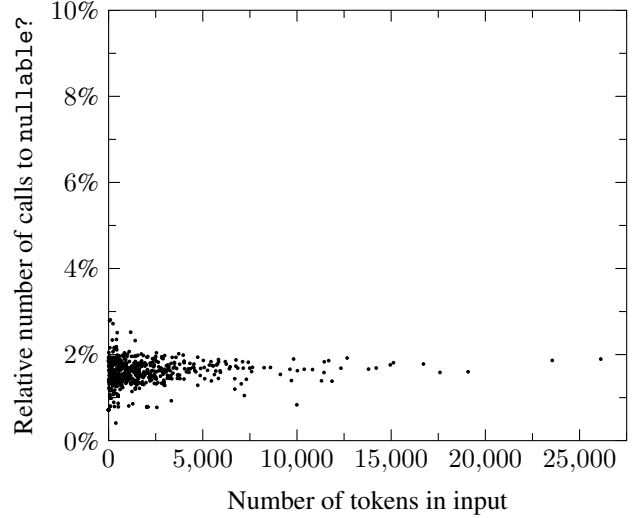


Figure 7. Number of calls to `nullable?` in the improved PWD relative to the original PWD

for Kleene algebras, however. They must preserve the structure of the resulting parse tree, and several of the identities insert reductions (\hookrightarrow) to do this.

4.3.1 Right-hand Children of Sequence Nodes

In our implementation, the following two reductions, which are used by Might et al. (2011), are not used during parsing. We omit these reductions because the forms on their left-hand sides cannot occur during parsing unless the initial grammar contains them.

$$\begin{aligned} p \circ \epsilon_s &\Rightarrow p \hookrightarrow \lambda u. (u, s) \\ p \circ \emptyset &\Rightarrow \emptyset \end{aligned}$$

Theorem 10. *While parsing, grammar nodes are never of the form $p \circ \epsilon_s$ or $p \circ \emptyset$ unless nodes in the initial grammar are of that same form.*

Proof. The derivative process changes only the left-hand child of a sequence node. Thus, the right-hand child of a sequence node is always a copy of the right-hand child of a sequence node from the initial grammar. \square

We take advantage of this fact by using these reduction rules on the initial grammar before parsing so that once parsing starts, we never need to check for them again. This avoids the need to inspect the right-hand children of sequence nodes during parsing and saves us the cost of any resulting memory accesses or conditional branching.

4.3.2 Canonicalizing Chains of Sequence Nodes

Consider a grammar fragment, like in Figure 8a, where p_1 is not nullable. When taking the derivative, only the left-hand children of the sequence nodes are considered. Thus, none

of p_2, \dots, p_{i-1}, p_i are inspected by `derive`, though the sequence nodes containing them are traversed. We could avoid the cost of this traversal if we restructured the grammar like in Figure 8b where the f' function rearranges the pairs in the resulting parse tree to match the AST produced by Figure 8a. As a result, `derive` would traverse only two nodes, the reduction node and topmost sequence node, instead of the i nodes in Figure 8a.

We can use compaction to try to optimize Figure 8a into Figure 8b by adding the following reduction rule, which implements associativity for sequence nodes.

$$(p_1 \circ p_2) \circ p_3 \Rightarrow (p_1 \circ (p_2 \circ p_3))$$

$$\hookrightarrow \lambda u. \{((t_1, t_2), t_3) \mid (t_1, (t_2, t_3)) \in u\}$$

However, this is not enough on its own. Depending on the order in which nodes get optimized by this reduction rule, a reduction node may be placed between neighboring sequence nodes that interferes with further applications of this reduction rule. This can lead to structures like in Figure 9a. Indeed, our inspection of intermediate grammars during parses revealed several examples of this.

We resolve this by also adding the following rule that floats reduction nodes above and out of the way of sequence nodes.

$$(p_1 \hookrightarrow f) \circ p_2 \Rightarrow$$

$$(p_1 \circ p_2) \hookrightarrow \lambda u. \{(f \{t_1\}, t_2) \mid (t_1, t_2) \in u\}$$

If we apply this rule for $(p_1 \hookrightarrow f) \circ p_2$ to the reduction nodes generated by applying the rule for $(p_1 \circ p_2) \circ p_3$, then we get Figure 9b where each f'_i does the work of f_i at the appropriate point in the AST. If we further use the rule for $(p \hookrightarrow f) \hookrightarrow g$ on the stack of reduction nodes in Figure 9b, we get Figure 8b, which allows derivatives to be computed efficiently.

Note that there is also a version of this reduction rule for when a reduction node is the right-hand instead of left-hand child of a sequence. It is the following.

$$p_1 \circ (p_2 \hookrightarrow f) \Rightarrow$$

$$(p_1 \circ p_2) \hookrightarrow \lambda u. \{(t_1, f \{t_2\}) \mid (t_1, t_2) \in u\}$$

However, for the same reasons as in Section 4.3.1, we use this only on the initial grammar and not during parsing.

4.3.3 Avoiding Separate Passes

Might et al. (2011) implement compaction as a separate pass in between the calls to `derive` for successive tokens. However, this means that nodes are traversed twice per token instead of only once. To avoid this overhead, we immediately compact nodes as they are constructed by `derive`. This results in two complications.

The first complication is that we do not want to iterate these reductions to reach a fixed point. We just do the reductions locally on the grammar node being generated by

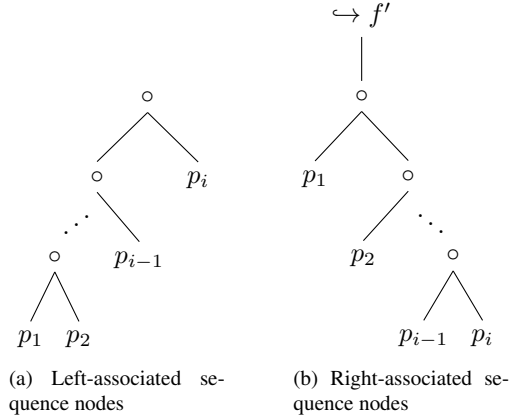


Figure 8. Examples of stacked sequence nodes

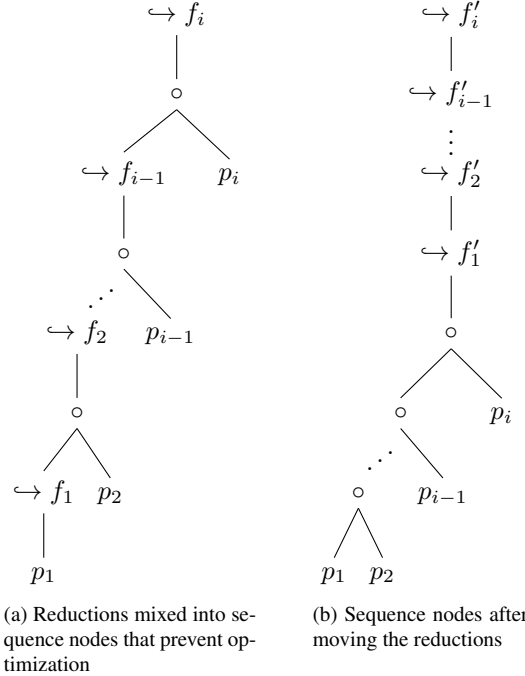


Figure 9. Examples of reductions mixed with sequence nodes

`derive`. As a result, there may be a few missed opportunities for applying reductions, but compactions in later derivatives should handle these.

The second complication is that we must consider how to compact when `derive` has followed a cycle in the grammar. The `derive` function usually does not need to know anything about the derivatives of the child nodes, which means that calculating these derivatives can be deferred using the lazy techniques described in Section 2.5. This poses a problem with compaction though, as many of the rules require knowing the structure of the child nodes. Like with the first complication, we have `derive` punt on this issue. If inspecting a child would result in a cycle, `derive` does not attempt

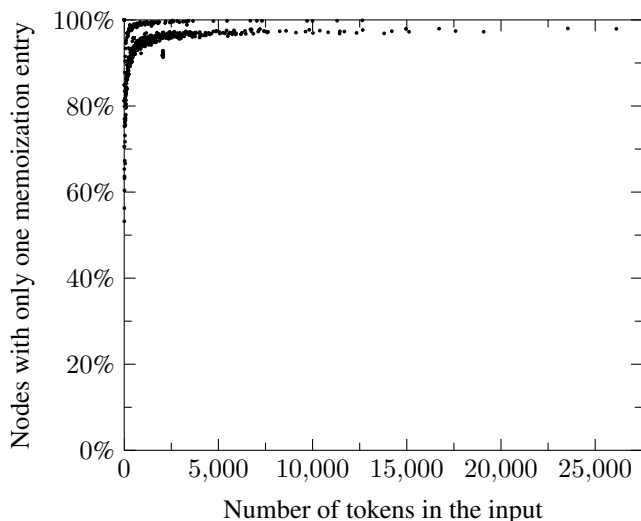


Figure 10. Percentage of nodes with only one memoization entry for derive

to compact. This design may miss opportunities to compact, but it allows us to avoid the cost of a double traversal of the grammar nodes.

4.4 Hash Tables and Memoization

The implementation in Might et al. (2011) uses hash tables to memoize `nullable?`, `derive`, and `parse-null`. Function arguments are looked up in those hash tables to see if a result has already been computed and, if so, what that result is. Unfortunately, hash tables can be slow relative to other operations. For example, in simple micro-benchmarks we found that that Racket’s implementation of hash tables can be up to 30 times slower than field access. Since memoization-table access is so central to the memoization process, we want to avoid this overhead. We do so by storing memoized results as fields in the nodes for which they apply instead of in hash tables mapping nodes to memoized results.

This technique works for `nullable?` and `parse-null`, but `derive` has a complication. The `derive` function is memoized over not only the input grammar node but also the token by which that node is being derived. Thus, for each grammar node, there may be multiple memoized results for multiple different tokens. The implementation used by Might et al. (2011) handles this using nested hash tables. The outer hash table maps grammar nodes to inner hash tables that then map tokens to the memoized result of `derive`. While we can eliminate the outer hash table by storing the inner hash tables for `derive` in a field in individual grammar nodes, the central importance of `derive` makes eliminating both sorts of hash table desirable.

These inner hash tables are usually small and often contain only a single entry. Figure 10 shows the percentage of inner hash tables in the original PWD implementation that

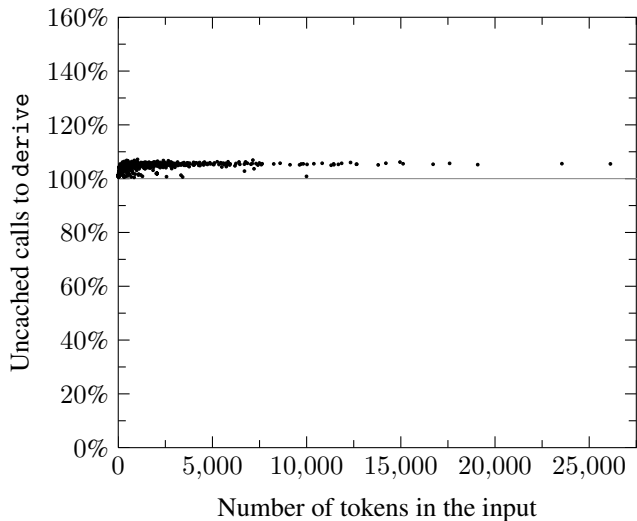


Figure 11. Percentage of uncached calls to derive with single entry versus full hash table

have only a single entry when parsing files from the Python Standard Library. Though we note the grouping into two populations, what interests us is that so many have only one entry. We can optimize for the single-entry case by adding two fields to each grammar node that behave like the key and value of a hash table that can store only one entry, and when a second entry is added, evicts the old entry.

This makes our memoization forgetful, and it may fail to notice when a token is reused multiple times in the input. However, the complexity results in Section 3 still hold, as they already assume every token is unique. Cycles in the grammar still require that we not forget the memoizations of `derive` on the current input token, but that requires only the single entry we store in each node.

We discovered that, in practice, the number of extra calls to `derive` that are recomputed as a result of this is relatively small. Figure 11 shows the number of calls to `derive` in our implementation when using the single-entry technique relative to the number when using full hash tables. While there are more uncached calls when using the single-entry technique, the increase is on average only 4.2% and never more than 4.8%. We also experimented with larger caches (e.g., double- or triple-entry caches) to see if the extra cache hits outweighed the extra computation cost. Early results were not promising, however, so we abandoned them in favor of a single-entry cache.

We measured the performance impact of this by running our implementation both with the single-entry technique and with full hash tables. The relative speedup of using the single-entry technique is shown in Figure 12. The extra calls to `derive` partially cancel out the performance improvements from avoiding the inner hash tables, but on average the performance still speeds up by a factor of 2.04.

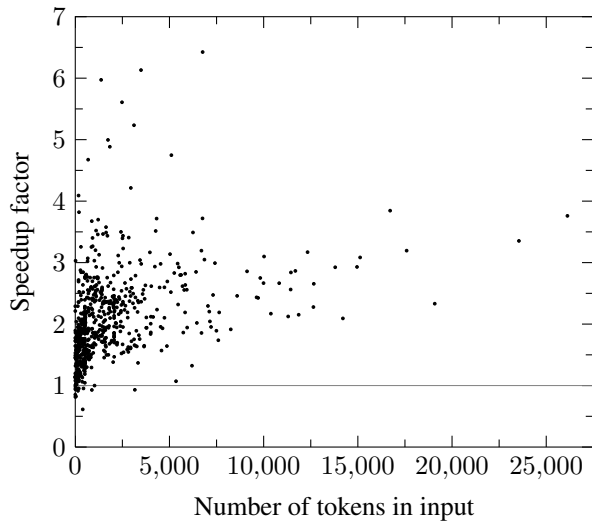


Figure 12. Performance speedup of single entry over full hash tables

5. Conclusion

In this paper, we have discovered that the believed poor performance of PWD both in theory and practice is not inherent to PWD. Rather, its worst-case performance at $O(n^3)$ is comparable to other full CFG parsers. Furthermore, with only a few algorithmic tweaks, the unacceptably slow performance of the implementation in Might et al. (2011) can be sped up by a factor of 951 to be on par with other parsing frameworks.

Acknowledgments

This material is partially based on research sponsored by DARPA under agreements number AFRL FA8750-15-2-0092 and FA8750-12-2-0106 and by NSF under CAREER grant 1350344. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon.

References

- Bison. Bison. URL <https://www.gnu.org/software/bison/>.
- Janusz A. Brzozowski. Derivatives of regular expressions. *Journal of the ACM (JACM)*, 11(4):481–494, October 1964. ISSN 0004-5411. doi: 10.1145/321239.321249.
- William Byrd. relational-parsing-with-derivatives, 2013. URL <https://github.com/webyrd/relational-parsing-with-derivatives>.
- Russ Cox. Yacc is not dead. Blog, December 2010. URL <http://research.swtch.com/yaccalive>.
- Jay Earley. *An Efficient Context-Free Parsing Algorithm*. PhD thesis, Carnegie Mellon University, 1968. URL <http://reports-archive.adm.cs.cmu.edu/anon/anon/usr/ftp/scan/CMU-CS-68-earley.pdf>.
- Jay Earley. An efficient context-free parsing algorithm. *Communications of the ACM*, 13(2):94–102, February 1970. ISSN 0001-0782. doi: 10.1145/362007.362035.
- Mark Engelberg. instaparse, 2015. URL <https://github.com/Engelberg/instaparse>.
- Abraham Flaxman, Aram W. Harrow, and Gregory B. Sorkin. Strings with maximally many distinct subsequences and substrings. *The Electronic Journal of Combinatorics*, 11(1):R8, 2004. ISSN 1077-8926. URL <http://www.combinatorics.org/ojs/index.php/eljc/article/view/v11i1r8>.
- Bryan Ford. Parsing expression grammars: a recognition-based syntactic foundation. In *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '04*, pages 111–122, New York, NY, USA, January 2004. ACM. ISBN 1-58113-729-X. doi: 10.1145/964001.964011.
- Gary A. Kildall. A unified approach to global program optimization. In *Proceedings of the 1st Annual ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages, POPL '73*, pages 194–206, New York, NY, USA, October 1973. ACM. doi: 10.1145/512927.512945.
- Dexter Kozen. A completeness theorem for kleene algebras and the algebra of regular events. *Information and Computation*, 110(2): 366–390, May 1994. ISSN 0890-5401. doi: 10.1006/inco.1994.1037.
- Bernard Lang. Deterministic techniques for efficient non-deterministic parsers. In Prof. Dr.-Ing. J. Loeckx, editor, *Automata, Languages and Programming*, volume 14 of *Lecture Notes in Computer Science*, pages 255–269. Springer Berlin Heidelberg, 1974. ISBN 978-3-540-06841-9. doi: 10.1007/3-540-06841-4_65.
- Tommy McGuire. Java-Parser-Derivatives, 2012. URL <https://github.com/tmmcguire/Java-Parser-Derivatives>.
- Gary H. Merrill. Parsing non- $LR(k)$ grammars with yacc. *Software: Practice and Experience*, 23(8):829–850, August 1993. ISSN 1097-024X. doi: 10.1002/spe.4380230803.
- Matthew Might. derp documentation, 2013. URL <http://matt.might.net/teaching/compilers/spring-2013/derp.html>.
- Matthew Might, David Darais, and Daniel Spiewak. Parsing with derivatives: a functional pearl. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming, ICFP '11*, pages 189–195, New York, NY, USA, September 2011. ACM. ISBN 978-1-4503-0865-6. doi: 10.1145/2034773.2034801.
- Russell Mull. parsing-with-derivatives, 2013. URL <https://github.com/mullr/parsing-with-derivatives>.
- Scott Owens, John Reppy, and Aaron Turon. Regular-expression derivatives re-examined. *Journal of Functional Programming*, 19(02):173–190, March 2009. ISSN 1469-7653. doi: 10.1017/S0956796808007090.
- Parser Tools. Parser tools: Context-free parsers. URL <http://docs.racket-lang.org/parser-tools/Context-Free-Parsers.html>.

Gregory Pfiel. *YACC-is-dead*, 2015. URL <https://github.com/sellout/YACC-is-dead>.

Python Software Foundation. Python 3.4.3, 2015a. URL <https://www.python.org/downloads/release/python-343/>.

Python Software Foundation. The Python language reference: Full grammar specification, 2015b. URL <https://docs.python.org/3/reference/grammar.html>.

Racket. Racket: Download racket v6.1.1, November 2014. URL <http://download.racket-lang.org/racket-v6.1.1.html>.

Frank Shearar. *Parsing-Derivatives*, 2013. URL <https://github.com/frankshearar/Parsing-Derivatives>.

Daniel Spiewak. Personal correspondence, October 2011.

Per Vognsen. *parser*, 2012. URL <https://gist.github.com/pervognsen/815b208b86066f6d7a00>.