

# Pika Parser: A parallelizable bottom-up packrat parser employing dynamic programming

LUKE A. D. HUTCHISON, Massachusetts Institute of Technology (alum)

A recursive descent or top-down parser is built from a set of mutually-recursive functions where each function directly implements one of the nonterminals of a grammar. The code structure of a recursive descent parser exactly parallels the grammar. Recursive descent parsers, or their memoized version, packrat parsers, are extremely simple to write, but suffer from some significant problems, in particular that left-recursive grammars cause the parser to get stuck in infinite recursion, and if there is a syntax error partway through a parse, it is difficult or impossible to optimally recover the parsing state to allow the parser to continue parsing past the error. We present the pika parser, a new type of parser that solves these problems by using dynamic programming to work bottom-up from character terminals towards the root of the parse tree. We describe a new longest-match PEG operator, which allows rules to be written directly in left-recursive form (simplifying grammar design), due to a new insight relating precedence order directly to match length. Bottom-up parsing conveys near-perfect error recovery capabilities on the parser, since all parsing depends only on the immediate surrounding context in the input. A pika parser has performance equivalent within a constant factor of a packrat parser, and is therefore linear in the size of the input and the depth of the parse tree. Since pika parsing employs dynamic programming (DP), parsing can be highly parallelized across the DP wavefront. Numerous new insights are presented into precedence, associativity, and left recursion in PEG grammars.

CCS Concepts: • **Theory of computation** → **Grammars and context-free languages**.

Additional Key Words and Phrases: parsing, packrat parsing, recursive descent parsing, top-down parsing, bottom-up parsing, precedence, associativity, left-recursive grammars, parsing error recovery, memoization, dynamic programming, parallel computing

## ACM Reference Format:

Luke A. D. Hutchison. 2020. Pika Parser: A parallelizable bottom-up packrat parser employing dynamic programming. *ACM Trans. Program. Lang. Syst.* TBD, TBD, Article TBD (TBD 2020), 30 pages. <https://doi.org/TBD>

## 1 INTRODUCTION

### 1.1 Parsing generative grammars

Since the earliest beginnings of computer science, computer scientists have sought to identify robust algorithms for parsing formal languages [1, 2]. Careful work to establish the theory of computation [24] provided further understanding into the expressive power of different classes of languages, the relationships between different classes of languages, and the relative computational power of algorithms that can parse different classes of languages.

Generative systems of grammars, particularly Chomsky’s context-free grammars (CFGs) [4] and Kleene’s regular expressions [16] became the workhorses of parsing of languages and data formats, and have been used nearly ubiquitously over several decades. However, the theory of these languages and the algorithms used to parse them are relatively complex, and additionally, inherent ambiguities and nondeterminism increase the implementational complexity of

---

Author’s address: Luke A. D. Hutchison, [luke.hutch@alum.mit.edu](mailto:luke.hutch@alum.mit.edu), Massachusetts Institute of Technology (alum).

---

Permission to make digital or hard copies of all or part 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 the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

parsers for these classes of language [20, 25], sometimes resulting in parsing conflicts, or exponential blowup in time or space [16, 21].

## 1.2 Bottom-up parsing

Much early work in language recognition focused on *bottom-up parsing*, in which as the input is processed in order, a forest of parse trees is built upwards from terminals towards the root, and these trees are joined together into larger trees as higher-level grammar rule *productions* or *reductions* are able to be applied. This parsing pattern is termed *shift-reduce* [1]. Examples of shift-reduce parsers include early work on precedence parsers and operator precedence parsers [3, 17, 23], the later development of LR [18] and LALR (lookahead) [10] parsers, and the more complex GLR parser [20] designed to handle ambiguous and nondeterministic grammars. Some parsers in this class run into issues with shift-reduce conflicts and/or reduce-reduce conflicts caused by ambiguous grammars; the GLR parser requires careful and complex data structure maintenance to avoid exponential blowup in memory usage, and to keep parsing time complexity to  $O(n^3)$  in the length of the input. Error recovery in all shift-reduce parsers is complex and difficult to implement optimally.

A particularly interesting bottom-up parser is the CYK algorithm [6, 7, 15] for context-free grammars. This parsing algorithm employs dynamic programming to efficiently recursively subdivide the input string into all substrings, and determines whether each substring can be generated from each grammar rule. This results in  $O(n^3|G|)$  performance in the length of the input  $n$  and the size of the grammar  $|G|$ . The algorithm is robust and simple, and in theory has near-perfect error recovery capability since all rules are matched against all substrings of the input. However, the CYK algorithm is inefficient for very large input sizes  $n$ , and additionally, the algorithm requires the grammar to be reduced to Chomsky Normal Form (CNF) [5] as a preprocessing step, which can result in a large or even an exponential blowup in grammar size in the worst case – so in practice,  $|G|$  may be large.

## 1.3 Top-down parsing

More recently, *recursive descent parsing* or *top-down parsing* has increased in popularity, due to its simplicity of implementation, and due to the nice property that the structure of mutually-recursive functions of the parser directly parallels the structure of the corresponding grammar. However, recursive descent parsers suffer from three significant problems:

- (1) The worst case time complexity of a recursive descent parser can be exponential in the length of the input.
- (2) Left-recursive grammars cause naive implementations of recursive descent parsing to get stuck in infinite recursion, which eventually causes a stack overflow.
- (3) If there is a syntax error partway through a parse, it is difficult or impossible to optimally recover the parsing state to allow the parser to continue parsing past the error. This makes recursive descent parsing particularly problematic for use in Integrated Development Environments (IDEs), where syntax highlighting and code assist need to continue to work between and after any syntax errors that are encountered during parsing.

These problems have gradually been addressed over the years:

- (1) The discovery of *packrat parsing*, employing memoization to avoid duplication of work [22], reduced parsing time for recursive descent parsing from potentially exponential in the length of the input to linear in the length of the input and the depth of the parse tree.

- (2) Left recursive rules were found to be rewritable as *indirectly* left recursive [11], although the rewriting rules can be extremely complex, especially when semantic actions are involved – or as right-recursive, with only a weak equivalence to the left recursive form. With extension, packrat parsing can be made to handle left recursion directly [11, 13, 14, 26], although usually with loss of linear time parsing, which is one of the most important reasons packrat parsing is chosen over other parsing algorithms [26].
- (3) Progress has been made on error recovery for recursive descent parsers, by identifying synchronization points [9], or by employing heuristic algorithms to attempt to recover after an error (yielding an “acceptable recovery rate” from 41% to 76% [8] in recent work). This is far from a solved problem.

## 1.4 PEG grammars

A significant refinement of recursive descent parsing known as Parsing Expression Grammars (PEG) was proposed by Ford in his 2002 PhD thesis [11]. PEG rules replace ambiguity with deterministic choice, i.e. all PEG grammars are unambiguous. There are only a handful of PEG rule types (see Section 1.4), and they can be implemented in a straightforward way. All non-left-recursive PEG grammars can be parsed directly by recursive descent, and they can be parsed efficiently with a packrat parser, i.e. using memoization. PEG grammars are more powerful than regular expressions, but it is conjectured (and as yet unproved) that there exist context-free languages that cannot be recognized by a PEG parser [12]). As PEG grammars are a subset of all grammars that can be parsed with recursive descent, PEG packrat parsers suffer from the same limitations as other packrat parsers: they cannot parse left-recursive grammars without special handling, and they have poor error recovery properties.

In contrast with most bottom-up parsers, which represent grammar rules as generative productions, PEG parsers represent grammar rules as *greedily-recognized patterns*. A PEG grammar consists of:

- A finite set  $N$  of *rule names*.
- A finite set  $\Sigma$  of *terminal symbols* that is disjoint from  $N$ . These could match individual characters of the input ('+'), words ("import"), or even a more complex terminal pattern recognized by a regular expression or similar, creating a hybrid parser.
- A *starting expression*  $e_s$ . This corresponds to the root node of the Abstract Syntax Tree (AST) after a successful parse.
- A finite set  $P$  of *parsing rules*. Each rule in  $P$  has the form  $A \leftarrow e$ , where  $A$  is a unique rule name and  $e$  is a *parsing expression or clause*.

A parsing expression or clause is a hierarchical expression consisting of:

- Terminal symbols.
- Rule names (as a reference to the rule of that name) – this corresponds to a recursive call in top-down parsing.
- The *empty string*  $\epsilon$  (with ASCII notation `()`). This always matches, even at the end of input, and consumes zero characters.
- a hierarchy of *PEG operators*, with parameters selected from the three items listed above. Each PEG operator and its parameters forms a *clause*. The parameters of the PEG operator are termed *subclauses*. Parentheses may be used to control precedence of nested operators.

Given clauses  $e$  and  $e_i$ , a new clause can be constructed using the following PEG operators:

Name	Num subclauses	Notation	Equivalent to	Example rule in ASCII notation
Seq	2+	$e_1 e_2 e_3$		Sum <- Prod '+' Prod;
First	2+	$e_1 / e_2 / e_3$		ArithOp <- '+' / '-';
Optional	1	$e?$	$e/\epsilon$	Value <- Name ('[' Num ']')?;
OneOrMore	1	$e+$		Whitespace <- [ \t\n\r ]+;
ZeroOrMore	1	$e^*$	$(e+)?$ <b>[or]</b> $e+\epsilon$	Program <- Statement*;
FollowedBy	1	$\&e$		Exclamation <- Word '&!';
NotFollowedBy	1	$!e$		NonExclamation <- Word '!';
Longest	2+	$e_1   e_2   e_3$		Expr <- E1   E2   E3;

Optional and ZeroOrMore can either be implemented directly, or automatically transformed into the equivalent forms shown above.

The match logic for these PEG operators is as follows:

- A Seq clause matches if all of its subclauses match the input, starting at the current parsing position. Each matching subclause consumes zero or more characters, which determines the starting parse position for the match attempt of the subsequent subclause. Matching stops once a single clause fails to match the input at its start position.
- A First clause matches if any of its subclauses match, in left to right order, with all match attempts starting at the current parsing position. After the first match is found, attempts to match further subclauses are abandoned. If no subclauses match, First clause does not match.
- An Optional clause matches if its subclause matches. If the subclause does not match, the Optional clause still matches (i.e. this operator always matches), but the match consumes zero characters.
- A OneOrMore clause matches if its subclause matches at least once, greedily consuming as many matches as possible.
- A ZeroOrMore clause matches if its subclause matches at least once, greedily consuming as many matches as possible, but if its subclause does not match, the ZeroOrMore still matches the input (i.e. this operator always matches), but consumes zero characters.
- A FollowedBy clause matches if its subclause matches at the current parsing position, but it does not consume any characters if it matches.
- A NotFollowedBy clause matches if its subclause does not match at the current parsing position, but it does not consume any characters if it matches.
- Longest is a new PEG operator type we propose (Section 1.5). A Longest clause matches if any of its subclauses match at the current parsing position. If multiple subclauses match, the subclause that consumes the longest span of input is used as the match, tiebreaking in left-to right order.

### 1.5 The new Longest PEG operator

Longest is a new PEG operator we propose as a way to fix left recursion issues and to simplify grammar-writing. The match of a Longest clause is simply the longest match of all of its subclauses.

Top-down recursion cannot handle directly left-recursive grammars, and even though bottom-up recursion (Section 2) has no problem with left recursion, the greedy nature of PEG operators causes parsing to fail for left-recursive grammars

on some inputs (Section 3.2). The Longest PEG operator solves this problem, while allowing grammars to still be written in left-recursive form (Section 3.4).

Additionally, without this new Longest operator, it can be difficult to get a standard PEG grammar to match the input as intended, particularly because the greedy nature of `First` can cause parsing failures if its subclauses were not listed in the correct order by the grammar-writer. An earlier subclause of a `First` clause may inadvertently match a prefix of the input that was intended to be matched by a later subclause, because the earlier clause matches prefixes of the later clause. To resolve this, the grammar writer must understand the partial ordering induced by all possible prefix matches between all subclauses of the `First` clause across all valid inputs, and then this partial ordering must be serialized into a total ordering using a topological sort, so that subclauses that match longer sequences are listed before subclauses that match prefixes of those longer sequences. If it is not easy or possible to determine this partial ordering of prefix matches of all subclauses across all inputs, then a Longest operator may be able to be used as a shortcut to obtain the desired parsing behavior. Longest rather than `First` will tend to match the grammar writer’s intent and intuition, since the PEG philosophy is one of greedily matching as many characters in the input as possible with a grammar rule.

## 2 THE PIKA PARSER: BOTTOM-UP PACKRAT PARSING

In this work we introduce the *pika parser*<sup>1</sup>, a new type of packrat parser that works bottom-up using dynamic programming, rather than top-down using basic memoization as with standard packrat parsing.

### 2.1 Advantages of pika parsing

Pika parsers are able to solve each of the predominant problems of other parsers:

- (1) As a near-semantically-equivalent “inversion” of top-down packrat parsing, pika parsing has roughly the same performance as packrat parsing, i.e. is roughly linear in the size of the input and the depth of the parse tree, usually within a moderate constant factor.
- (2) Bottom-up parsing directly handles left-recursive grammars and grammars containing arbitrary numbers of cycles, which can greatly simplify grammar design.
- (3) Bottom-up parsing conveys near-perfect error recovery capabilities on the parser, since all parsing depends only on the immediate surrounding context in the input. In fact it is straightforward to recover parsing at any level of the grammar, by skipping ahead from a syntax error to the next valid match of a grammar rule of interest. This will be particularly helpful in providing code assist features for IDEs that work robustly in spite of syntax errors.
- (4) We define the pika parser to work with PEG grammars, which are unambiguous and deterministic, eliminating ambiguity and nondeterminism problems that affect CFG parsers.
- (5) The dynamic programming algorithm employed by the pika parser enables parsing to be parallelized across the input, or more accurately, to be parallelized across a dynamic programming “wavefront” which initially consists of all terminals in the input. This allows multiple parts of the input to be parsed simultaneously, splitting the work across all available cores for most of the parsing process. With the right data structures, this parallelization is very straightforward.

---

<sup>1</sup>A pika (/ˈpaɪkə/) is a small mountain-dwelling rabbit-like mammal found in Asia and North America, typically inhabiting moraine fields near the treeline. Like a packrat, it stores food in caches or “haystacks” during the warmer months to survive the winter.

## 2.2 High-level overview of pika parsing

The premise of pika parsing is very simple: take the recursive structure of a packrat parser's parse of a given input with a given grammar, and invert the top-down recursive structure so that parsing happens bottom-up from terminals towards the root. Instead of shallower frames of recursion recursing down into deeper frames of recursion, as with recursive descent parsing, bottom-up parsing has deeper frames of recursion trigger the execution of shallower frames of recursion. This seems like a strange concept, until one has the realization that this is exactly the definition of dynamic programming.

In most dynamic programming (DP) algorithms, there is an underlying recurrence, which is represented as lookups into earlier, already-computed parts of the DP table. These are directly analogous to recursive calls in recursive descent parsing. As the DP algorithm proceeds, it follows these dependency arrows backwards from leaves (usually the top and left edge of the table, if the recurrence is two-dimensional) towards the root (represented by the bottom right cell of the table). Because the recurrence is the same in every cell, figuring out how to reverse these dependency arrows is very simple – but the problem is further simplified since all DP table cells are always evaluated, so you only have to find an evaluation order so that all the previous values needed for an evaluation of the recurrence relation are always available. For a DP algorithm like the string edit distance algorithm [19], the recurrence relation for cell  $\langle i, j \rangle$  depends upon cells  $\langle i - 1, j \rangle$ ,  $\langle i, j - 1 \rangle$  and  $\langle i - 1, j - 1 \rangle$ , so it is possible to populate the DP table in row major order, column major order, or using a trailing diagonal "wavefront" that sweeps from upper left to lower right through the table. Any of these orders ensures that a "lower" cell in the recurrence is always evaluated before a "higher" cell in the recurrence.

Inversion of packrat parsing applies exactly the same idea as dynamic programming, except that the structure of the recurrence relations (i.e. the recursive calls) between the cells are now dynamic rather than static. For a Seq clause like  $P \leftarrow A B C$ , the starting position at which subclause B will be matched depends upon the number of characters consumed by the match of subclause A, and the starting position of subclause C depends upon how many characters are consumed by both A and B. These dynamic relationships are recorded by storing backRefs or back-references in the memo table entries, which record when a parent clause tried to match a subclause at a given position.

In other words, to produce a bottom-up "inverted recursion" when there is a dynamic dependency relationship between parent clauses and subclauses, each step in the bottom-up parse actually performs a single top-down step from recursion depth  $i$  to recursion depth  $i + 1$ , in order to check the contents of the memo table for subclause matches at a deeper level of recursion. Then back-references are recorded from that deeper level of recursion, so that the parent clause can be found in the future if the subclause match is later improved.

Bottom-up parsing is probably best explained with a concrete example.

## 2.3 Example of pika parsing

Listing 1 shows a simple "sum of products" grammar. The most complex expression this could recognize is a sum of two products. The lowest precedence rule S (i.e. Sum) is listed as the first (root) rule, so the root of the parse tree will be an S node if the parse is successful. The S rule is comprised of two P rules (i.e. Product), which are higher precedence. The P rule contains two clauses: an outer First clause and (in the first position) a nested Seq clause,  $(N \ X \ N)$ . The N rule recognizes one or more digits. The A and X rules recognize arithmetic operators.

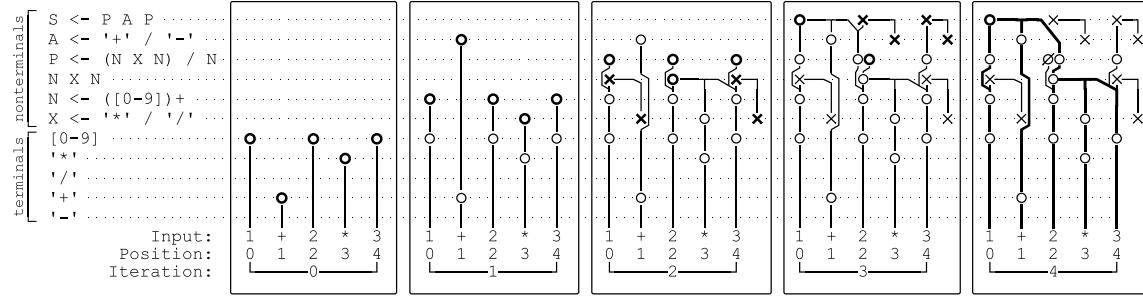


Fig. 1. Example parse of Listing 1

```

1 Grammar:
2   S <- P A P;
3   A <- '+' / '-';
4   P <- (N X N) / N;
5   N <- ([0-9])
6   X <- '*' / '/';
7 Input:
8   1+2*3

```

Listing 1. Example grammar and input string

Iteration	activeSet entry	Match(o)/Mismatch(x)	Match len	Parent memo keys
0	$\langle [0-9], 0 \rangle$	o	1	$\langle N, 0 \rangle$
0	$\langle '+', 1 \rangle$	o	1	$\langle A, 1 \rangle$
0	$\langle [0-9], 2 \rangle$	o	1	$\langle N, 2 \rangle$
0	$\langle '*', 3 \rangle$	o	1	$\langle X, 3 \rangle$
0	$\langle [0-9], 4 \rangle$	o	1	$\langle N, 4 \rangle$
1	$\langle N, 0 \rangle$	o	1	$\langle N \ X \ N, 0 \rangle, \langle P, 0 \rangle$
1	$\langle A, 1 \rangle$	o	1	
1	$\langle N, 2 \rangle$	o	1	$\langle N \ X \ N, 2 \rangle, \langle P, 2 \rangle$
1	$\langle X, 3 \rangle$	o	1	
1	$\langle N, 4 \rangle$	o	1	$\langle N \ X \ N, 4 \rangle, \langle P, 4 \rangle$
2	$\langle N \ X \ N, 0 \rangle$	x	–	
2	$\langle P, 0 \rangle$	o	1	$\langle S, 0 \rangle$
2	$\langle N \ X \ N, 2 \rangle$	o	3	$\langle P, 2 \rangle$
2	$\langle P, 2 \rangle$	o	1	$\langle S, 2 \rangle$
2	$\langle N \ X \ N, 4 \rangle$	x	–	
2	$\langle P, 4 \rangle$	o	1	$\langle S, 4 \rangle$
3	$\langle S, 0 \rangle$	o	3	
3	$\langle P, 2 \rangle$	o	3	$\langle S, 0 \rangle$
3	$\langle S, 2 \rangle$	x	–	
3	$\langle S, 4 \rangle$	x	–	
4	$\langle S, 0 \rangle$	o	5	

Table 1. Memo table entries created in an example pika parse of Listing 1.

Fig. 1 and Table 1 show how pika parsing proceeds with this example. Each of the five boxes in Fig. 1 depicts the content of the memo table at the end of one step of parsing. The grammar rules are shown on the left hand side. Note that the rule  $P$  contains the nested clause  $(N \times N)$ , the clause  $N$  contains the nested terminal  $[0-9]$ , and the clauses  $A$  and  $X$  contain nested terminals for the arithmetic operators, and each of these nested clauses are each given their own row in the memo table (this is so that duplicate work is avoided in more complex grammars, where sub-clauses may be shared between multiple rules). In Fig. 1, matches in the memo table are shown by the symbol  $\circ$  and mismatches are shown by the symbol  $\times$  (note that the use of  $\times$  to indicate mismatch is unrelated to the rule with the name  $X$ ). These symbols are bolded if they are new matches added to the memo table at the end of the iteration.

A *memo key* or index into the memo table has the form  $\langle \text{clause}, p \rangle$ , where *clause* designates the row of the memo table and *p* designates the column (or the start position within the input of a match of the given clause).

Between iteration  $i$  and iteration  $i+1$ , any new matches or mismatches added to the memo table (marked in bold) must trigger the evaluation of all applicable *parent memo keys* to be evaluated in the next iteration. This is the mechanism by which the top-down PEG parsing algorithm is inverted and turned into a bottom-up algorithm. These parent memo keys are added to a set, *activeSet*, which is the upwards-expanding wavefront of the dynamic programming algorithm (also known sometimes as the “dirty set” when describing dynamic programming algorithms). Since *activeSet* elements are independent of each other (they all depend only on the state of the memo table at the end of the previous iteration, not upon each other), *activeSet* elements can all be processed in parallel.

The first mechanism for identifying parent memo keys is termed *seedParentClauses*. Given the rules  $S \leftarrow P \ A \ P$  and  $P \leftarrow (N \ X \ N) / N$ , rule  $S$  is a seed parent clause of rule  $P$ , since in top-down evaluation, an attempt to match rule  $S$  at position 0 would first try to match subclause  $P$  at position 0. Therefore, in bottom-up parsing, if the memo table entry with the memo key  $\langle P, p \rangle$  is updated in iteration  $i$  for some position  $p$ , then  $\langle S, p \rangle$  must be scheduled for matching in iteration  $i+1$ . The entries in *seedParentClauses* are determined for each clause and subclause at initialization, since the grammar does not change during parsing (this *seedParentClauses* is static). The parent memo key consists of the seed parent clause of the subclause, paired with the start position  $p$  of the subclause match’s memo key.

The second mechanism for identifying parent memo keys is termed *backRefs*, and is dynamic, in the sense that the start position of the parent memo key (stored in *backRefs* within a memo entry) depends upon how many characters the parent clause consumed before it tried to match the subclause in question. When any parent clause attempts to match the input at a given start position, for each subclause that the parent clause tries to match by looking up the subclause at the current parsing position in the memo table, a back reference memo key is stored in the subclause’s memo entry at that position, linking back to the parent clause and the parent clause’s start position (whether or not the subclause actually matched at that start position). If the subclause match changes in a subsequent iteration (due to identifying a better match for the subclause at this start position), then any back-referenced parent clauses are triggered for re-evaluation.

The update steps for the memo table are described below.

**2.3.1 Iteration 0.** The memo table is initially seeded with all matches of each terminal at some position of the input, requiring  $O(N|T|)$  time for input length  $N$ , and given  $|T|$  terminals. (Optimizations of this preprocessing step will be discussed later.) This yields the following memo table entries. Note that parent memo keys produced from each *activeSet* entry in iteration  $i$  become the *activeSet* entries for iteration  $i+1$ .

**2.3.2 Iteration 1.** Each parent memo key in the *activeSet* again matches a single character at a higher clause in the grammar, triggering their own parent clauses for evaluation in Iteration 2. However,  $\langle A, 1 \rangle$  does not have any



seedParentClauses, because it is the second clause of  $S \leftarrow P \ A \ P$ , and the first clause consumes at least one character (similarly  $\langle X, 3 \rangle$  does not have any seedParentClauses).

**2.3.3 Iteration 2.** Each of the new active set entries  $\langle P, 0 \rangle$ ,  $\langle P, 2 \rangle$  and  $\langle P, 4 \rangle$  generated by iteration 1 results in an attempt to match  $P$  at the corresponding position. Since  $P$  is a `First` clause, this results in first testing whether the first subclause of  $P \leftarrow (N \ X \ N) / N$ , i.e.  $N \ X \ N$ , matches at position  $p \in 0, 2, 4$ . As always, this is achieved by looking at *the state of the memo table at the end of the previous iteration*. The memo table has no entries for  $N \ X \ N$  at the end of iteration 1, so all three of these subclause match attempts fail, and then the second subclause  $N$  of rule  $P$  is tested. In all three cases, there is a positive match entry in the memo table at the end of iteration 1; therefore,  $\langle P, p \rangle$  is recorded as a match in the memo table for all three values of  $p$  at the end of iteration 2, and the parent memo entries  $\langle S, p \rangle$  are triggered by adding them to `activeSet` for the next iteration.

However, in spite of the fact that  $\langle N \ X \ N, p \rangle$  does not match in the memo table for any of  $p \in 0, 2, 4$  at the end of iteration 1, when the clause  $\langle P, p \rangle$  looks up those subclause memo entries, a `backRef` is added to each subclause memo entry, pointing back to  $\langle P, p \rangle$ . If the match state of any of the  $\langle N \ X \ N, p \rangle$  entries subsequently changes from a mismatch to a match, or from a shorter match to a longer match, the corresponding parent clause  $\langle P, p \rangle$  will be triggered in the following iteration.

For `activeSet` entries  $\langle N \ X \ N, 0 \rangle$  and  $\langle N \ X \ N, 4 \rangle$ , the first subclause  $N$  matches (as determined by looking up  $\langle N, 0 \rangle$  and  $\langle N, 4 \rangle$  respectively in the memo table in the state it was in at the end of iteration 1). In both cases, the subclause match consumes one character, so when the second subclause  $X$  is matched, it starts matching at the following character position (1 and 5 respectively). However,  $\langle X, 1 \rangle$  does not match, and  $\langle N, 5 \rangle$  is past the end of the input. Entries are added to `backRefs` for these mismatching entries, pointing back to their respective parent clauses (in this case, the subclause match status will never change, so these back references will never be followed).

A complete match is found for `activeSet` entry  $\langle N \ X \ N, 2 \rangle$ , based on three entries in the memo table at the end of iteration 1:  $\langle N, 2 \rangle$ ,  $\langle X, 3 \rangle$ , and  $\langle N, 4 \rangle$ , consuming a total of three characters, matching the last three characters of the input string, "2\*3". This match causes the parent memo key  $\langle P, 2 \rangle$  to be added to `activeSet` for the next iteration.

**2.3.4 Iteration 3.** The `activeSet` entries  $\langle S, 2 \rangle$  and  $\langle S, 4 \rangle$  fail to match, because they fail to find a match of rule  $A$  in their second subclause (input position 3 and 5 respectively). However,  $\langle S, 0 \rangle$  succeeds in matching the first three characters, i.e.  $1 + 2$ , via the entries in the memo table at the end of iteration 2:  $\langle P, 0 \rangle$ ,  $\langle A, 1 \rangle$ , and  $\langle P, 2 \rangle$ . Note that this match does not respect precedence order, and consequently it will be discarded in the next iteration and replaced with a better match. Rule  $S$  does not have any parent clauses, so this match does not cause any new memo keys to be written to `activeSet`.

The `activeSet` entry  $\langle P, 2 \rangle$ , i.e. the rule  $P \leftarrow (N \ X \ N) / N$ , checks the memo table for a match of its first subclause,  $\langle N \ X \ N, 2 \rangle$ , and finds a match in that position, consuming three characters. Therefore,  $\langle P, 2 \rangle$  is a new match in iteration 3, also consuming 3 characters.

Note that  $\langle P, 2 \rangle$  matched only one character at the end of iteration 2 (the previous best match), but now matches 3 characters at the end of iteration 3 (the new best match). This is shown as two adjacent `o` characters in the memo table at this position – the older match is shown to the left, non-bolded, and the newer match is shown to the right, bolded. Because memo table values are always looked up in the state the memo table was in at the end of the previous iteration, the match for  $\langle S, 0 \rangle$  actually depends upon the older match obtained in iteration 2, not upon the new best match found in the iteration 3. The previous best match will be overwritten by the new best batch at the end of iteration 3 (shown in iteration 4 by putting a slash through the older match). Replacing the best match in this memo entry will trigger the

parent clause  $\langle S, \emptyset \rangle$  to be re-evaluated in iteration 4, via the backRefs mechanism – hence this is parent memo key generated by  $\langle P, 2 \rangle$  in this iteration.

**2.3.5 Iteration 4.** The single  $\langle S, \emptyset \rangle$  is re-evaluated, and now for its third subclause, picks up the updated match for  $\langle P, 2 \rangle$  from the memo table in the state it was in at the end of iteration 3. The final subclause matches are  $P = "1"$ ,  $A = "+"$ , and  $P = "2*3"$ , corresponding to the parse  $"1+(2*3)"$ . The final parse tree is shown with bolded lines.

At this point activeSet is empty, because the clause of the only active memo entry  $\langle S, \emptyset \rangle$  has no parent clauses, so the parse is complete.

**2.3.6 Efficiency Analysis.** The final parse tree, constructed bottom-up, has the same structure as the parse tree that would have been obtained by top-down parsing of the same input with the same grammar<sup>2</sup>.

There are a number of entries in the memo table at the end of the parse that are unused (meaning they are not part of the final parse tree rooted at  $\langle S, \emptyset \rangle$ ), such as mismatching entries and entries that were eagerly triggered by bottom-up parsing that did not end up in the final parse tree. Two entries in the memo table,  $\langle S, \emptyset \rangle$  and  $\langle P, 2 \rangle$ , started with a worse match, and the match was replaced in a later iteration with a better match. (The logic used to determine whether a match is better or worse is given later.) In each of these cases, unnecessary work was done by the bottom-up pika parser, compared to a top-down packrat parser.

The overhead of wasted work performed by bottom-up parsing is not enormous – it should amount to only a moderately small constant factor. This is the price paid for the advantages listed in Section 2.1, and the performance overhead of bottom-up parsing can be mitigated by parallelization of the parser.

Not all the unused nodes left over in the memo table at the end of parsing are useless: in particular, in the case of an incomplete parse due to a syntax error, all matches for any desired grammar rule can be found by simply reading off the entries in the row of the memo table corresponding to the rule. Thus an IDE or compiler can very easily recover after a syntax error at any desired level of grammar, for example by finding the start of the next valid expression, statement, or function.

Additionally, the replacing of an earlier shorter match with a later longer match in the same memo entry is not always wasteful of work, because the earlier shorter match may be incorporated into the parse tree as a descendant of the later longer match, as in the case of a grammar that contains cycles, such as an expression rule that can match of a self-recursive tree of sub-expressions of mutually nested types. The matches stored in the memo table are actually connected to their subclause matches, i.e. they form nodes in the cumulatively built structure of the parse tree.

**2.3.7 Preprocessing optimization: lexing.** The preprocessing step of initializing the memo table with all matches of all terminals against the input can be optimized by running a lex (or lexing) preprocessing step, which is a shallow acyclic top-down preprocessing parse of the input stream to break it into complete tokens. This can reduce the number of memo entries that the memo table is seeded with.

Lexing serves another important purpose: the PEG operator OneOrMore greedily matches from left to right, consuming as many matches of its single subclause as possible. Given the rules  $\text{Word} \leftarrow [A-Z]^+$  and  $\text{Num} \leftarrow [0-9]^+$  that matches a sequence of one or more letters or digits respectively, and given the input CAT123, the bottom-up parsing algorithm would start by seeding the memo table with the matches for the terminal clauses  $[A-Z]$  and  $[0-9]$ , i.e.  $\{\langle 'C', 0 \rangle, \langle 'A', 1 \rangle, \langle 'T', 2 \rangle, \langle '1', 3 \rangle, \langle '2', 4 \rangle, \langle '3', 5 \rangle\}$ . Each of the letter matches would seed the  $[0-9]^+$  parent clause and each

<sup>2</sup>This is not quite true, especially considering that left-recursive grammars can be parsed bottom-up without any problem, but top-down recursion of the same grammar can get stuck in an infinite loop. However, assuming a top-down parse succeeds, the parse tree obtained by bottom-up parse should be structurally the same as the parse tree obtained by the top-down parse.

of the digit matches would seed the  $[0-9]^+$  parent clause. Therefore, the first step of bottom-up parsing would produce the memo table matches  $\{ \langle \text{"CAT"}, 0 \rangle, \langle \text{"AT"}, 1 \rangle, \langle \text{'T'}, 2 \rangle, \langle \text{"123"}, 3 \rangle, \langle \text{"23"}, 4 \rangle, \langle \text{'3'}, 5 \rangle \}$ , in other words all suffix strings of every `OneOrMore` token is added to the memo table. This can cascade up the memo table until all the suffixes are finally discarded by a larger rule combining matches from multiple adjacent rules, e.g. `WordNum <- Whitespace? WordNum Whitespace?`, which selects only the first of each series of suffix matches.

An even more important use of lexing is to deal with quoted strings and comments. Given that bottom-up parsing starts with no context, a language consisting of tokens, quoted strings and comments may incorrectly assume that a word inside a quoted string or comment is a token. In general, the contents of a quoted string or comment can produce numerous unnecessary memo table entries in the lower levels of parsing. These unnecessary entries will eventually be ignored as the structure of the document starts to emerge at higher and higher levels of parsing, as more and more adjacent context is joined. However, it is wasteful to store these unnecessary matches in the memo table. If the lex rule extracts tokens, quoted strings and comments in a preprocessing parse, and if these lexed tokens are used to seed the bottom-up parse rather than character-level terminals, this problem is solved.

In the reference parser (Section 5), lexing is enabled by adding a rule called `Lex` that matches all terminals (or tokens, quoted strings, comments, etc.) using a single `First` clause. This rule is repeatedly applied until all input is consumed. If a `Lex` rule is used, all terminals should be included somewhere in the `Lex` rule, so that the memo table is properly seeded with any terminal match that should end up in the final parse tree, otherwise the appropriate parent clauses will not be triggered for subsequent bottom-up parsing.

Note that top-down parsing of the `Lex` rule in the reference parser implements non-memoized recursive descent parsing (to avoid creating unnecessary memo entries, in order to keep the memo table small), so the tokens matched by the `Lex` rule should be short, and each clause of the `Lex` rule should be simple, in order to keep lex preprocessing running in close to linear time.

## 2.4 Implementation details

The key classes of the reference pika parser are as follows<sup>3</sup>. The code shown is simplified for to illustrate the most important aspects of the parsing algorithm and data structures. The full code for this class can be viewed in the reference parser implementation (see Section 5).

**2.4.1 The MemoKey class.** This class is a simple record type holding the two-dimensional coordinate of a memo table entry, i.e.  $\langle \text{clause}, \text{startPos} \rangle$ .

---

```

1  class MemoKey {
2      Clause clause;
3      int startPos;
4  }
```

---

Listing 2. The MemoKey class

**2.4.2 The Parser class.** This class accepts a `Grammar` in its constructor (not shown). A `Grammar` is a collection of `Rule` objects, and also has an `allClauses` field, which is a set of all unique clauses and sub-clauses across all rules. The

<sup>3</sup>Listings are given in Java-like notation, however some Java-specific details and quirks have been simplified or omitted, for example all data structures that may be written from multiple threads are prefixed by `Concurrent` for clarity, even if a class of this exact name does not exist in the Java standard libraries. For example, `ConcurrentPriorityQueue` represents the Java standard library class `java.lang.concurrent.PriorityBlockingQueue`, and `ConcurrentSet` represents a concurrent hashset resulting from calling `Collections.newSetFromMap(new ConcurrentHashMap<T, Boolean>())`. Some field initializers (new statements) and modifiers such as `public` and `final` have been omitted, as well as trivial assignment constructors.

main parse method of the `Parser` class first creates two sets, `activeSet` (which holds the memo keys constituting the current DP wavefront), and `updatedEntries`, which holds the set of memo entries for which new matches were found in the current iteration, then initializes the memo table, by populating the memo table with the locations that all terminals match the input. (Note that the optional lex preprocessing step is not shown here, since it is only an optimization.)

(Note the use of the Java construction `collection.parallelStream().forEach(lambda)` here, which processes all elements of `collection` in parallel by calling the function `lambda` on each element. Both the initialization step and each stage of the main parsing loop can be item-wise parallelized using this construction, since none of the items in the collections being iterated over have any interdependencies.)

After initialization, the parse method enters the main parsing loop. This loop consists of two alternately-run inner loops (each implemented as a parallel `forEach`) that are mutually interdependent upon each other's output. The first loop iterates through and consumes all items in `activeSet`, the current DP wavefront, checking all subclauses of each clause to see if they match in the memo table, and applying the PEG operator for each clause to determine if the clause itself matches based on the subclause match(es). This first loop produces `updatedEntries`, the set of all memo table entries for which a new match was found.

The second inner loop iterates through and consumes all `updatedEntries`, overwriting the previous best match for each updated memo table entry with the new best match, generating the next `activeSet` from the parent clauses of these new best matches. Overwriting the previous best match with the new best match is performed in this second separate step to ensure the memo table does not change while memo entries are still being read for the `activeSet` entries in the first inner loop.

Each element of `activeSet` generates elements in `updatedEntries`, and vice versa, until both sets are empty, at which point the parse is complete, because there are no more parent clauses to trigger, and no more memo table entries to update.

---

```

1  public class Parser {
2      final Grammar grammar;
3      MemoTable memoTable = new MemoTable();
4
5      public void parse(String input) {
6          var activeSet = new ConcurrentSet<MemoKey>();
7          var updatedEntries = new ConcurrentSet<MemoEntry>();
8
9          // Initialize memo table by matching all terminal clauses against the input at each position.
10         // Adds the memo entries for the initial terminal matches to updatedEntries.
11         grammar.allClauses.parallelStream().forEach(clause -> {
12             if (clause instanceof Terminal) {
13                 for (int startPos = 0; startPos < input.length(); startPos++) {
14                     var memoKey = new MemoKey(clause, startPos);
15                     // Add a match to the memo table if the terminal matches
16                     clause.match(MatchDirection.TOP_DOWN, memoTable, memoKey, input, updatedEntries);
17                 }
18             }
19         });
20
21         // Main parsing loop
22         while (!activeSet.isEmpty() || !updatedEntries.isEmpty()) {
23             // Test match status of each memo key in the active set, producing updatedEntries
24             activeSet.parallelStream().forEach(memoKey -> {
25                 memoKey.clause.match(MatchDirection.BOTTOM_UP, memoTable, memoKey, input,
26                                     updatedEntries);
27             });
28             activeSet.clear();
29
30             // Update the best match for all updatedEntries, populating activeSet with parent clauses
31             updatedEntries.parallelStream().forEach(memoEntry -> {
32                 memoEntry.updateBestMatch(input, activeSet);
33             });
34             updatedEntries.clear();
35         }
36     }
37 }

```

---

Listing 3. The Parser class

**2.4.3 The Clause class.** This class represents a PEG operator or clause and its subClauses. This class is abstract, and is subclassed for each PEG operator type (the subclasses implement the match method).

The `Clause.seedParentClauses` field records which parent clauses try to match this subclause at the parent clause’s start position (i.e. before the parent clause has consumed any previous characters). Like `MemoEntry.backRefs`, `Clause.seedParentClauses` triggers parent clauses to be checked for a match when the match for one of their sub-clauses is updated. These two back-reference mechanisms are collectively used to turn top-down recursion into bottom-up recursion. Both are needed, since `MemoEntry.backRefs` stores back-references to parent clauses no matter how many previous characters the parent clause previously consumed, whereas `Clause.seedParentClauses` is needed to trigger the evaluation of parent clauses that have not yet been evaluated for a match (and therefore have not stored back references in their subclauses’ memo entries).

The `canMatchZeroChars` field is used to record whether a clause can match while consuming zero characters, such as  $\epsilon$ ,  $X?$ ,  $X^*$ , or  $X / Y / Z?$ . (Clauses that can match zero characters always match at every input position, consuming zero or more characters, which slightly changes how they should be handled, as will be described later.)

The `match` method is used to look up a subclause match in the memo table if `matchDirection == BOTTOM_UP`, or to do top-down recursive descent parsing for lexing if `matchDirection == TOP_DOWN`.

---

```

1  abstract class Clause {
2      Clause[] subClauses;
3      Set<Clause> seedParentClauses;
4      boolean canMatchZeroChars;
5
6      abstract Match match(MemoTable memoTable, MemoKey memoKey, String input,
7                          Set<MemoEntry> updatedEntries);
8  }

```

---

Listing 4. The Clause class

**2.4.4 The Char class.** This class is the simplest example of a PEG terminal clause. `Char` extends `Terminal`, and `Terminal` extends `Clause`. The `Terminal` class is used to detect which clauses are terminals using the `instanceof` operator.

The `Char` class simply attempts to match the character it is initialized with at the location provided by `memoKey.startPos`. If the character matches, it stores a one-character match in the memo table at the entry indexed by `memoKey`, by calling `memoTable.addTerminalMatch`.

---

```

1  class Char extends Terminal {
2      char c;
3
4      Char(char c) {
5          this.c = c;
6      }
7
8      @Override
9      Match match(MatchDirection matchDirection, MemoTable memoTable, MemoKey memoKey, String input,
10                 Set<MemoEntry> updatedEntries) {
11          return memoKey.startPos < input.length() && input.charAt(memoKey.startPos) == c
12              ? memoTable.addTerminalMatch(memoKey, /* terminalLen = */ 1, updatedEntries)
13              : null;
14      }
15  }

```

---

Listing 5. The implementation of the Char terminal

**2.4.5 Non-terminal PEG operator classes.** The `First` subclass of `Clause` implements the `First` PEG operator. The `match` method iterates through each subclause of the `First` clause, and returns a new `Match` object for the first subclause that matches, or `null` if no match was found. In contrast, the `Seq` subclass' `match` method requires all subclauses to match in order for the `Seq` clause to match, and the match start position must be updated for each subclause as characters are consumed.

If matching top-down (i.e. during lexing), the `match` method recursively calls the `match` method of the subclause. When matching bottom-up, the `match` method looks for subclause matches in the memo table by calling `memoTable.lookupMemo` method, but does not recurse to the subclauses' `match` methods.

The other subclasses of Clause (implementing the other PEG operator types) can be found in the reference implementation (see Section 5).

---

```

1  class First extends Clause {
2      @Override
3      Match match(MatchDirection matchDirection, MemoTable memoTable, MemoKey memoKey,
4                  String input, Set<MemoEntry> updatedEntries) {
5          for (int subClauseIdx = 0; subClauseIdx < subClauses.length; subClauseIdx++) {
6              var subClause = subClauses[subClauseIdx];
7              var subClauseMemoKey = new MemoKey(subClause, memoKey.startPos);
8              var subClauseMatch = matchDirection == MatchDirection.TOP_DOWN
9                  ? subClause.match(MatchDirection.TOP_DOWN, memoTable, subClauseMemoKey,
10                                     input, updatedEntries)
11                  : memoTable.lookupMemo(subClauseMemoKey, input, memoKey, updatedEntries);
12              if (subClauseMatch != null) {
13                  return memoTable.addNonTerminalMatch(memoKey,
14                                                         /* firstMatchingSubClauseIdx = */ subClauseIdx,
15                                                         new Match[] { subClauseMatch }, updatedEntries);
16              }
17          }
18          // No match
19          return null;
20      }
21  }

```

---

Listing 6. The implementation of the First PEG operator

---

```

1  class Seq extends Clause {
2      @Override
3      Match match(MatchDirection matchDirection, MemoTable memoTable, MemoKey memoKey, String input,
4                  Set<MemoEntry> updatedEntries) {
5          var subClauseMatches = new Match[subClauses.length];
6          var currStartPos = memoKey.startPos;
7          for (int subClauseIdx = 0; subClauseIdx < subClauses.length; subClauseIdx++) {
8              var subClause = subClauses[subClauseIdx];
9              var subClauseMemoKey = new MemoKey(subClause, currStartPos);
10             var subClauseMatch = matchDirection == MatchDirection.TOP_DOWN
11                 ? subClause.match(MatchDirection.TOP_DOWN, memoTable, subClauseMemoKey,
12                                    input, updatedEntries)
13                 : memoTable.lookupMemo(subClauseMemoKey, input, memoKey, updatedEntries);
14             if (subClauseMatch == null) {
15                 // If a subclause fails to match, the Seq clause does not match
16                 return null;
17             }
18             subClauseMatches[subClauseIdx] = subClauseMatch;
19             currStartPos += subClauseMatch.len;
20         }
21         return memoTable.addNonTerminalMatch(memoKey, /* firstMatchingSubClauseIdx = */ 0,
22                                                subClauseMatches, updatedEntries);
23     }
24 }

```

---

Listing 7. The implementation of the Seq PEG operator

**2.4.6 The MemoEntry class.** This class stores the current best match, and a priority queue of new best matches that were found in the current iteration. The `addNewBestMatch` method adds new matches during the first inner loop of the main parse loop, when `activeSet` is being iterated over. Importantly, if a new match is not better than the current best match (i.e. if `newMatch.compareTo(bestMatch) >= 0`, see the explanation of the `Match.compareTo` method below), then the new match is never written to the `newMatches` priority queue. This ensures that every new best match *monotonically* improves the match, so that the parse is guaranteed to terminate.

The `updateBestMatch` method is called by the second inner loop of the main parse loop, when `updatedEntries` is being iterated over. The `updateBestMatch` method finds the best new match, by removing the head of the `newMatches` priority queue, then setting `bestMatch` to this new best match. Then all the seed parent clauses of the clause of the new best match and all the backrefs from the memo entry are used to generate parent memo keys that are added to the `activeSet` for the next iteration.



---

```

1  class MemoEntry {
2      MemoKey memoKey;
3      Match bestMatch;
4      ConcurrentPriorityQueue<Match> newMatches = new ConcurrentPriorityQueue<>();
5      ConcurrentSet<MemoKey> backRefs = new ConcurrentSet<>();
6
7      MemoEntry(MemoKey memoKey) {
8          this.memoKey = memoKey;
9      }
10
11     void addNewBestMatch(Match newMatch, Set<MemoEntry> updatedEntries) {
12         // If the new match is better than the current best match from the previous iteration
13         if (bestMatch == null || newMatch.compareTo(bestMatch) < 0) {
14             // Add the new match to the priority queue, and mark the entry as updated
15             newMatches.add(newMatch);
16             updatedEntries.add(this);
17         }
18     }
19
20     void updateBestMatch(String input, Set<MemoKey> activeSetOut) {
21         // Get the best new updated match for this MemoEntry, if there is one
22         var newBestMatch = newMatches.peek();
23         if (newBestMatch != null) {
24             newMatches.clear();
25             // Replace previous best match with new best match
26             bestMatch = newBestMatch;
27             // Generate new activeSet entries from seed parent clauses and backrefs
28             for (var seedParentClause : memoKey.clause.seedParentClauses) {
29                 MemoKey parentMemoKey = new MemoKey(seedParentClause, bestMatch.memoKey.startPos);
30                 activeSetOut.add(parentMemoKey);
31             }
32             for (var backref : backRefs) {
33                 activeSetOut.add(backref);
34             }
35             // Clear backrefs
36             backRefs.clear();
37         }
38     }
39 }

```

---

Listing 8. The MemoEntry class

**2.4.7 The Match class.** This class is used to store matches and their lengths in memo entries. Match instances are nodes in the parse tree, connected to their child nodes by their subClauseMatches field. The most important method in Match is `compareTo`, which is used to determine when one match is better than another. This comparator method is used by the `MemoEntry.newMatches` priority queue so that the best new match is always at the head of the queue, and also to ensure the monotonic improvements of memo entry updates, as described above.

According to the Comparable contract, the comparison `match1.compareTo(match2)` should return a negative number if `match1` should sort earlier than `match2` (in this context, if `match1` is a better match than `match2`, zero if the two matches are equal, and a positive number if `match1` should sort later than `match2` (in this context, if `match1` is a worse match than `match2`). This `compareTo` method is designed to as closely as possible replicate the greedy left-to-right semantics of the PEG operators as they are defined for top-down parsing.

---

```

1  class Match implements Comparable<Match> {
2      MemoKey memoKey;
3      int len;
4      Match[] subClauseMatches;
5      int firstMatchingSubClauseIdx;
6
7      @Override
8      public int compareTo(Match o) {
9          if (o == this) {
10             return 0;
11         }
12         // An earlier subclause match in a First clause wins over a later match
13         if (this.memoKey.clause instanceof First) {
14             var diff0 = this.firstMatchingSubClauseIdx - o.firstMatchingSubClauseIdx;
15             if (diff0 != 0) {
16                 return diff0;
17             }
18         }
19         // A longer match wins over a shorter match
20         var diff1 = o.len - this.len;
21         if (diff1 != 0) {
22             return diff1;
23         }
24         var minSubclauses = Math.min(this.subClauseMatches.length, o.subClauseMatches.length);
25         for (int i = 0; i < minSubclauses; i++) {
26             var diff2 = o.subClauseMatches[i].len - this.subClauseMatches[i].len;
27             if (diff2 != 0) {
28                 return diff2;
29             }
30         }
31         // A longer list of OneOrMore subclause matches wins over a shorter list
32         return o.subClauseMatches.length - this.subClauseMatches.length;
33     }
34 }

```

---

Listing 9. The Match class

**2.4.8 The MemoTable class.** This class contains a double-level map, `memoTable`, mapping a clause to a start position to a `MemoEntry`. This nesting of maps, and the use of `ConcurrentSkipListMap` for the inner map implementation, allows all matches anywhere in the input to be quickly enumerated for any clause, in order of start position. This is useful for error recovery and other purposes.

The `getOrCreateMemoEntry` method returns the memo entry for a given memo key, or if a memo entry doesn't yet exist for that memo key, creates a new memo entry for that memo key and returns the new memo entry. (The implementation shown below uses `computeIfAbsent`, which is the Java-specific mechanism for creating Map entries in a lock-free way while avoiding calling `new` to create a new entry where possible. It is guaranteed to implement singleton allocation semantics in this case.)

The `lookupMemo` entry simply looks up a memo key in the memo table, and returns the current best match (for the state the memo table was in at the end of the previous iteration). However, it also adds the memo key of the parent clause (the clause that called this method to check its subclause matches) to the `backRefs` field of the subclause's memo

entry. This dynamically provides for the parent clause to be triggered for re-evaluation if the memo entry's best match is ever updated.

The `addMatch` method (with variants for nonterminals and terminals) creates a new `Match` object, linking a match to its subclause matches, and then calls `MemoEntry.addNewBestMatch` to add the new match to the priority queue of new matches generated in the current iteration. The best of these matches will overwrite the memo entry's current `bestMatch` value at the end of the iteration, if it is an improvement over the previous best match.

---

```

1  class MemoTable {
2      Map<Clause, ConcurrentSkipListMap<Integer, MemoEntry>> memoTable = new ConcurrentMap<>();
3
4      // Get the memo entry for memoKey if one is already present; if not, create and return a new entry
5      private MemoEntry getOrCreateMemoEntry(MemoKey memoKey) {
6          var skipList = memoTable.computeIfAbsent(memoKey.clause, c -> new ConcurrentSkipListMap<>());
7          var memoEntry = skipList.computeIfAbsent(memoKey.startPos, s -> new MemoEntry(memoKey));
8          return memoEntry;
9      }
10
11     Match lookupMemo(MemoKey memoKey, String input, MemoKey parentMemoKey,
12                     Set<MemoEntry> updatedEntries) {
13         var memoEntry = getOrCreateMemoEntry(memoKey);
14         memoEntry.backRefs.add(parentMemoKey);
15         return memoEntry.bestMatch;
16     }
17
18     private Match addMatch(MemoKey memoKey, int firstMatchingSubClauseIdx, int len,
19                           Match[] subClauseMatches, Set<MemoEntry> updatedEntries) {
20         var memoEntry = getOrCreateMemoEntry(memoKey);
21         var newMatch = new Match(memoKey, firstMatchingSubClauseIdx, len, subClauseMatches);
22         memoEntry.addNewBestMatch(newMatch, updatedEntries);
23         return newMatch;
24     }
25
26     Match addNonTerminalMatch(MemoKey memoKey, int firstMatchingSubClauseIdx,
27                              Match[] subClauseMatches, Set<MemoEntry> updatedEntries) {
28         var totSubClauseMatchLen = 0;
29         for (Match subClauseMatch : subClauseMatches) {
30             totSubClauseMatchLen += subClauseMatch.len;
31         }
32         return addMatch(memoKey, firstMatchingSubClauseIdx, totSubClauseMatchLen, subClauseMatches,
33                         updatedEntries);
34     }
35
36     Match addTerminalMatch(MemoKey memoKey, int terminalLen, Set<MemoEntry> updatedEntries) {
37         return addMatch(memoKey, /* firstMatchingSubClauseIdx = */ 0, terminalLen, new Match[0],
38                         updatedEntries);
39     }
40 }

```

---

Listing 10. The `MemoTable` class

**2.4.9 Handling zero-length positive matches.** Actually the code shown above for `MemoTable.lookupMemo` is incomplete, because of an interesting corner case, involving positive matches that consume zero characters. The terminal  $\epsilon$  (or  $()$  in ASCII representation), which positively matches zero characters everywhere in the input (including beyond the end of

the input string) is the most basic example of this type of match. There are numerous other examples of PEG clauses that can positively match zero characters, such as `Optional`, since the clause `X?` is equivalent to `X / ()`, and if the clause `X` does not match, `()` will always match zero characters.

This can cascade up the grammar from terminals towards the root. For example, a `Seq` clause whose subclauses can all positively match zero characters can itself positively match zero characters. A `First` clause with even a single subclause that can positively match zero characters can itself positively match zero characters. (Furthermore, in the case of `First`, any subclauses after the first subclause that can positively match zero characters will be ignored, because that subclause will always match.)

The PEG operator that matches  $\epsilon$ , `Nothing`, is a terminal clause, so the memo table could be seeded with positive zero-length matches for  $\langle \text{Nothing}, p \rangle$  for all  $0 \leq p \leq (\text{input}.len)$ . This will lead to the grammar being parsed correctly bottom-up, since all parent clauses containing a `Nothing` terminal subclause will be properly triggered by seeding the memo table during initialization, just as with all other terminals. However, this adds a large number of unnecessary matches to the memo table, and these matches can cascade upwards through higher-level clauses that can also positively match zero characters, creating many unnecessary memo entries at higher levels of the grammar that will ultimately be ignored because most of them will not be connected to the final parse tree. Therefore, it is better to never memoize `Nothing` matches to avoid polluting the memo table.

Instead, we make the assumption that a rule will never have `Nothing` as its first subclause, which practically speaking is a reasonable assumption. This allows us to only schedule parent clauses for matching if `Nothing`, or any another clause that can positively match zero characters, is looked up in the memo table as the second or subsequent subclause of a clause. We do this by updating the code for `MemoTable.lookupMemo` as shown below. If a memo entry is looked up by a parent clause, and there is no known match for that entry, but the subclause can positively match zero characters, then a zero-length match is added to the memo table for the memo entry, by calling `memoEntry.addNewBestMatch`. This match will trigger the parent clause to be re-evaluated for matching in the next iteration.

Without this refinement, failing to memoize `Nothing` can cause parsing to fail, because the original definition of `lookupMemo` will return `null` (implying there is no match) for subclauses that can positively match zero characters but have never been memoized before (because they were not triggered bottom-up from the `Nothing` terminal).

This revised method requires the `canMatchZeroChars` field to be set for each clause during initialization. The DAG of all grammar clauses must be topologically ordered, then the initialization code should iterate upwards through the DAG of clauses from terminals towards the root grammar rule, determining which clauses can positively match zero characters.

---

```

1 Match lookupMemo(MemoKey memoKey, String input, MemoKey parentMemoKey,
2     Set<MemoEntry> updatedEntries) {
3     var memoEntry = getOrCreateMemoEntry(memoKey);
4     memoEntry.backRefs.add(parentMemoKey);
5     if (memoEntry.bestMatch != null) {
6         return memoEntry.bestMatch;
7     } else if (memoKey.clause.canMatchZeroChars) {
8         // Find index of first subclause that can match zero characters
9         int firstMatchingSubClauseIdx = 0;
10        for (int i = 0; i < memoKey.clause.subClauses.length; i++) {
11            if (memoKey.clause.subClauses[i].canMatchZeroChars) {
12                firstMatchingSubClauseIdx = i;
13                break;
14            }
15        }
16        // Create and memoize a new zero-width match. This will trigger the parent clause to be
17        // re-evaluated in the next iteration.
18        var newMatch = new Match(memoKey, firstMatchingSubClauseIdx, /* len = */ 0, new Match[0]);
19        memoEntry.addNewBestMatch(newMatch, updatedEntries);
20        return newMatch;
21    } else {
22        return null;
23    }
24 }

```

---

Listing 11. Improvement of the MemoTable.lookupMemo method to handle zero-length matches without memoizing Nothing

### 3 PRECEDENCE PARSING WITH PEG GRAMMARS

#### 3.1 Representing precedence levels with PEG rules

Because PEG parsers are naturally greedy, and operate in a top-down, left-to-right order in their standard recursive descent form, building a hierarchy of operator precedence into a grammar can be complicated. Because left recursive grammars result in infinite recursion, a grammar containing a left-recursive cycle must be rewritten in a non-left-recursive form, usually by using a OneOrMore operator, complicating the grammar, particularly if a parse tree contains more than one loop around the grammar cycle.

However, when PEG parsing is run bottom-up, cycles are handled naturally, so it becomes very natural to build precedence into the parser. Consider the following grammar, where  $E[i]$  defines grammar rule  $E$  for precedence level  $i$ , and  $E$  on the right hand side of the rule is a self-reference that respects operator precedence:

---

```

1 E[4] <- '(' E ')';
2 E[3] <- '[0-9]+' / '[a-z]+';
3 E[2] <- '-' E;
4 E[1] <- E ('*' / '/') E;
5 E[0] <- E ('+' / '-') E;

```

---

Listing 12. Example grammar containing cycles (self-references); rule names have precedence suffix

This grammar can be rewritten directly as standard PEG rules, where each reference to  $E$  is rewritten to refer to the next highest level of precedence of the grammar, as required by the semantics of precedence. If an entire underlying grammar rule fails, then the grammar simply defers to the next highest level of precedence, effectively elevating the match attempt to the next highest level of precedence.

(With top-down PEG parsing, this failover to higher precedence levels results in a depth-first postorder traversal of the grammatical structure of the input, so even though it is counterintuitive to start recursing at the lowest-precedence rule, this is the correct order to ensure proper precedence semantics. Elsewhere in the grammar, references to rule E should resolve to the lowest precedence clause, i.e. references to E from grammar rules other than E should rewrite E to  $E\emptyset$ , or should match the implicit rule  $E \leftarrow E\emptyset$ .)

---

```

1 E4 <- '(' E0 ')';
2 E3 <- '[0-9]+' / '[a-z]+' / E4;
3 E2 <- '-' E3 / E3;
4 E1 <- (E2 '*' / '/') E2 / E2;
5 E0 <- (E1 '+' / '-') E1 / E1;

```

---

Listing 13. Example grammar with precedence (does not handle associativity)

### 3.2 Representing both precedence and left-associativity with PEG rules

However, even though (with either top-down or bottom-up parsing) the above grammar can handle nested levels of precedence, e.g.  $1*2+3*4$ , it cannot handle runs of equal precedence, or direct nesting of operators of equal precedence, for example  $1+2+3$  or  $((2+3))$  or  $--3$ .

It is tempting to want to try to handle this situation by ensuring that one of the operands defers first to the same level of precedence, and then to the next highest level of precedence, for example it would be easy to assume that the following grammar implements left precedence for  $E\emptyset$  and  $E1$ :

---

```

1 E4 <- '(' (E4 / E0) ')';
2 E3 <- '[0-9]+' / '[a-z]+' / E4;
3 E2 <- '-' (E2 / E3) / E3;
4 E1 <- ((E1 / E2) '*' / '/') E2 / E2;
5 E0 <- ((E0 / E1) '+' / '-') E1 / E1;

```

---

Listing 14. Example grammar, handling associativity using First

This left-recursive grammar can match  $1+2+3$  as  $(1+2)+3$  as expected, as well as  $((2+3))$  and  $--3$ ; however, it cannot match the string  $1+2*(3+4)$ , because this requires multiple loops around the grammar cycle via the  $E4$  rule, and bottom-up parsing gets stuck before the second loop around the grammar can be matched. Note that a right-associative self-recursive precedence rule such as  $E\emptyset \leftarrow (E1 '+' (E\emptyset / E1)) / E1$  would not get stuck with bottom-up parsing; left recursion has special properties.

For the input  $1+2*(3+4)$ , the failure of bottom-up parsing due to left recursion occurs as follows:

- The initial terminal matches are added to the memo table, including  $\langle [0-9], 2 \rangle$  matching "2".
- $\langle [0-9]^+, 2 \rangle$  matches "2".
- $\langle E3, 2 \rangle$  matches "2".
- $\langle E2, 2 \rangle$  matches "2", through the final  $E3$  subclause of rule  $E2$ .
- $\langle E\emptyset, 5 \rangle$  matches  $3+4$ .
- $\langle E1, 2 \rangle$  fails to match  $2*(\dots)$ , because there is no  $\langle E4, 4 \rangle$  match yet. Therefore, the final "failover" subclause promotes the match to  $E2$ , and  $\langle E1, 2 \rangle$  is found to match "2".
- $\langle E4, 4 \rangle$  matches  $(3+4)$ .

- $\langle E0, 0 \rangle$  matches "1+2", because its third subclause  $\langle E1, 2 \rangle$  currently matches "2". This fragment of the putative parse tree does not respect precedence, so in theory the rightmost sub-clause match ("2") should be replaced with a better match in a future iteration, specifically with "2\*(3+4)". (That doesn't end up happening.)
- $\langle E1, 2 \rangle$  matches "2\*(3+4)". This is a new better match than the previous match for this memo entry, so it triggers  $\langle E0, 0 \rangle$  to be added to `activeSet` for re-matching in the following iteration, through the `backRefs` mechanism, since the third subclause match of  $\langle E0, 0 \rangle$  is  $\langle E1, 2 \rangle$ .
- When attempting to match  $\langle E0, 0 \rangle$  again,  $E0$ 's own directly left-recursive subclause ( $E0 / E1$ ) greedily trips up the match, because this subclause will now greedily consume the entire previous best match for  $\langle E0, 0 \rangle$  as its first operand, and will try to parse ('+' / '-' ) in the fourth character position, after "1+2"; therefore, the match fails.
- The  $E1$  rule cannot itself fix the problem by expanding its own best batch from "2\*(3+4)" to "1+2\*(3+4)", because the "1+" part of the input only matches at lower precedence.
- The memo table ends up with two locally-best matches that overlap by one character, "1+2" and "2\*(3+4)", and the parsing ends without producing a complete parse tree.

For similar reasons, this grammar cannot parse "1+2+3\*4\*5", due to an overlap of  $E0$  and  $E1$  matches at "3". (This overlap is somewhat analogous to a shift-reduce conflict in other bottom-up parsers, though has different properties.)

Furthermore, this grammar can parse directly-nested parentheses, e.g. " $((1+2))$ ", but cannot match left-recursive nested parentheses contained within a different precedence level match, e.g. " $((1+2)*3)$ ". This can be fixed by rewriting the highest-precedence rule as  $E4 \leftarrow '( (E4 \mid E0) )'$  (see Section 3.4).

There are two straightforward mechanisms we can use to avoid or fix the left-associativity problem without modifying the parser.

### 3.3 Rewriting left-recursive rules using right-recursive OneOrMore

The first solution to the left recursion problem is to simply avoid associativity altogether. We can take a clue from the rewriting of left-recursive rules into non-left-recursive (OneOrMore) form to guess at how the left recursion problem could be solved with grammar rewriting.

Directly applying standard heuristics used to rewrite left-recursive rules into non-left-recursive form would solve the left recursion problem not just for top-down parsing, but also for bottom-up parsing. However, these rewriting heuristics can be complex, *since they try to ensure that the produced parse tree maintains left-associative structure* – and for inputs that require multiple loops around grammar cycles, rewriting rules in this way can result in a blowup of the size and complexity of the grammar. There is a better alternative, if we relax the requirement that the parse tree should maintain left-associative structure.

As demonstrated, any rule at a given precedence level that has a left-recursive reference to its own precedence level will eventually fail, even with bottom-up parsing. The whole purpose of adding these self-references from a precedence level to itself is to add left associativity to the structure of the parse tree. However, associativity is only a convention used to specify how to interpret an arithmetic expression. Any notation that requires interpretation based on conventions is fundamentally ambiguous. A better way to approach this problem is to not even try to resolve the interpretation of equal-precedence runs of operands and operators until after parsing has completed – in other words, precedence should not even be resolved in the parse tree, it should be resolved if necessary in a postprocessing step,

perhaps when converting the parse tree into an Abstract Syntax Tree, or the Abstract Syntax Tree into a Concrete Syntax Tree.

In other words, clauses at a given precedence level should only ever reference higher precedence levels, never equal precedence levels – except in the case of a *precedence-breaking pattern* such as a rule that matches a set of parentheses around a lower-precedence clause.

Adjacent terms of equal precedence are then collected together into associativity-agnostic runs using a `OneOrMore` clause. For binary (two-arg) arithmetic expressions, these can be directly written as right-recursive rules, avoiding the complexity of the standard rule-rewriting heuristics, since we no longer care about resolving associativity as the parse tree is built:

---

```

1 E4 <- '(' (E4 | E0) ')';
2 E3 <- '[0-9]+' / '[a-z]+' / E4;
3 E2 <- ('-' E3) / E3;
4 E1 <- (E2 (('*' / '/') E2)+) / E2;
5 E0 <- (E1 (('+' / '-') E1)+) / E1;
```

---

Listing 15. Example grammar, handling associativity using right-recursive `OneOrMore`

Using this grammar, all adjacent terms of equal precedence (i.e. all `E0` and `E1` matches) will match as sibling subclauses of a `OneOrMore` match in the final parse tree.

This grammar will correctly parse `"1+2*(3+4)"` as `1+(2*(3+4))`, and also handles expressions containing runs of more than two terms of equal precedence, e.g. it will correctly parse `"1+2+3*4*5"` as `(1+2+(3*4*5))` (the parse tree in this case consists of ternary internal nodes rather than the usual binary nodes, due to the `OneOrMore` operators).

This grammar rewriting pattern is simple enough and general enough to satisfy most needs for precedence parsing without modifying the parser.

### 3.4 Rewriting left-recursive rules using `Longest`

There is a second mechanism that can be used to allow a pika parser to directly handle left recursion without any problem, using the newly proposed `Longest` PEG operator (Section 1.5).

Consider a hypothetical *nondeterministic recursive descent parser*, which is able to handle left recursion directly without resulting in infinite recursion. This parser would parse top-down, and would simultaneously follow both the left-recursive branch of recursion and the next non-left-recursive `First` subclause, e.g. Listing 14 (i.e. for grammar rule `E0`, if the first branch of the `(E0 / E1)` subclause fails to match, the match result of the second branch is used, although both branches are run simultaneously, and the `First` operator is only applied after both branches return). To ensure termination, and because nested precedence levels must monotonically decrease in length (due to Lemma 1), each successively nested level of left recursion must match one fewer character, starting at the top of the recursion with the span of the input string from the initial start position to the end of the input. The left recursion terminates when it matches zero characters.

The nondeterministic splitting of execution paths for this top-down algorithm would work directly with left recursive grammars, but would result in exponential execution time as a function of recursion depth on a deterministic Turing machine. Fortunately, when parsing bottom-up, the exponentially-large space of nondeterministic splits explored by top-down nondeterministic parsing can be directly pruned to a linear subset, because none of the failed branches of top-down nondeterministic parsing are ever seeded in the first place for bottom-up parsing.

We can now consider the relationship between match length and precedence.



DEFINITION 1. *If a top-down parser has recursed  $n$  times around a cycle of  $L$  precedence levels before reaching precedence level  $p$ , where  $0 \leq p < L$ , then the absolute precedence of the rule is  $p^* = -nL + p$ .*

(The negation of the first term in the absolute precedence formula is to ensure that deeper recursion frames, have lower precedence than shallower recursion frames, no matter the number of loops around the grammar of each.)

This definition admits the following lemma:

LEMMA 1. *Given two matches from the same precedence hierarchy, where the matches have the same starting position and different match lengths, the shorter match must be of higher absolute precedence than the longer match.*

Restated, the lengths of matches of different rules from a precedence hierarchy that share a start position are monotonically reverse-ordered with respect to the absolute precedence level of each match.

This lemma holds because each successively lower precedence level of the grammar adds its own terminals or subclause matches while incorporating higher-precedence matches, therefore a lower-precedence match consumes more input than a higher precedence match. The highest precedence rule loops the grammar using a precedence-breaking pattern like parentheses, which also consumes additional input, so the absolute precedence definition continues to hold across one or more loops around the grammar.

A consequence of Lemma 1 is that a left-recursive rule within the a precedence hierarchy that is intended for top-down parsing by a nondeterministic parser can be equivalently rewritten using Longest rather than First, for example the rule  $E0 \leftarrow ((E0 / E1) ('+' / '-')) E1 / E1$  can be equivalently rewritten as  $E0 \leftarrow ((E0 | E1) ('+' / '-')) E1 / E1$ .

The rewritten rule selects  $E0$  as the left operand if that match is longer than any match at the same start position for  $E1$ , which follows both precedence order and left associativity. If both clauses match, but the match for  $E1$  is longer than the match for  $E0$ , then the match for  $E1$  must have looped at least one more time around the precedence levels than the match for  $E0$ , i.e. the path through the parse tree passes through the precedence-breaking pattern (e.g. parentheses) of the highest precedence level at least one additional time for  $E1$  than for  $E0$ . Therefore,  $E1$  is actually lower precedence than  $E0$ , and  $E1$  should be used as the left subclause match, rather than  $E0$ .

Rewriting left-recursive rules in this way almost allows bottom-up parsing with a deterministic parser to succeed, although one further change is needed, because the greediness of PEG operators has opposite semantics in top-down vs. bottom-up.

The fundamental reason why a bottom-up parser was getting stuck with left-recursive subclauses (Section 3.2) is that directly comparing precedence levels, e.g.  $(E0 / E1)$ , only takes into account relative precedence level within the same loop around the grammar, not absolute precedence level. Rewriting just the leftmost subclause using Longest as shown above can still cause bottom-up parsing to fail in situations where greedily matching the longest left operand leads to a shorter match of the whole rule. The longest whole-rule match should be the one that succeeds, because it is the length of the whole rule at a given precedence level that is monotonically ordered with respect to the absolute precedence level, not the length of the subclauses. Therefore, we need an alternative rewriting heuristic to break make the greedy semantics of bottom-up parsing match the greedy semantics of top-down parsing.

Note that the above reasoning applies not just to directly left-recursive rules, but also *left-recursive structures in the input string*, when the matching rule has a single self-reference subclause. For example, as mentioned previously, the rule  $E4 \leftarrow '(' (E4 / E0) ')'$  cannot match left-recursive nested parentheses with a different precedence level match between them, e.g.  $((1+2)*3)$ : the rule is not left-recursive, but the structure of the input string is left-recursive. The reason this rule fails to match is that  $'(' (E4 / E0) ')'$  will greedily match  $'(' (E4$  against the first part of

the sequence, " $((1+2))$ ", then will fail to match the ')' that it expects. This can be fixed by rewriting this rule to  $E4 \leftarrow ('(E4 \mid E0)')$ , which chooses the longer match " $(1+2)*3$ " over the shorter match " $(1+2)$ ".

The correct rule rewriting pattern to enable left-recursive rules (and left-recursive inputs) to produce a left-associative parse tree using a bottom-up parser, is as follows. Given the rule  $A[i, L] \leftarrow A \text{ Op } A$ , where the extra "L" modifier denotes left associativity:

- Duplicate the rule using Longest; rewrite the duplicates into left-recursive and non-left-recursive form:  
 $A[i, L] \leftarrow (A[i] \text{ Op } A[i+1]) \mid (A[i+1] \text{ Op } A[i+1])$
- Add whole-rule failover to the next-highest precedence level using First:  
 $A[i, L] \leftarrow ((A[i] \text{ Op } A[i+1]) \mid (A[i+1] \text{ Op } A[i+1])) / A[i+1]$

(Note that  $A[i+1]$  here relies on the precedence levels wrapping around, in other words this next-highest-preference rule reference would be more accurately written as  $A[(i+1) \bmod N]$ , given  $N$  precedence levels ranging from 0 to  $(N - 1)$  inclusive – however, we omit this detail to simplify the notation.)

It is interesting to note that  $A[i, L] \leftarrow ((A[i] \text{ Op } A[i+1]) \mid (A[i+1] \text{ Op } A[i+1])) / A[i+1]$  does not work with Longest replacing First, i.e. in the form  $A[i, L] \leftarrow ((A[i] \text{ Op } A[i+1]) / (A[i+1] \text{ Op } A[i+1])) / A[i+1]$ , because in this form, the First construction only considers relative precedence levels *within the same loop around the precedence hierarchy*, whereas by Lemma 1, Longest compares absolute precedence levels by comparing match length, so the number of loops around the precedence hierarchy is taken into consideration. On the other hand, the failover clause  $A[i] \leftarrow \dots / A[i+1]$  does not need to take into account absolute precedence levels, except in the case of the highest precedence clause, since the purpose of this failover clause is to promote the match to the next highest absolute precedence level if the match fails, which is the same as the next highest relative precedence level for all but the highest precedence clause. This is why the highest precedence clause needs to use Longest rather than First to work with left-recursive inputs. It is faster to use First rather than Longest for the failover clause (and leaves fewer unused entries in the memo table), since First only matches as many subclauses as needed, but Longest always matches all its subclauses.

Rewriting a left-recursive rule to produce a right-associative parse tree is simpler. Given the rule  $A[i, R] \leftarrow A \text{ Op } A$ , where the extra "R" modifier denotes right associativity, the rewriting pattern is:

- Rewrite self-references to be right-recursive:  
 $A[i, R] \leftarrow A[i+1] \text{ Op } A[i]$
- Add whole-rule failover to the next-highest relative precedence level using First:  
 $A[i, R] \leftarrow (A[i+1] \text{ Op } A[i]) / A[i+1]$

For all but the highest precedence level, given a rule with a single self-referencing subclause, e.g.  $A[i] \leftarrow X \ A \ Y$ , the rewriting pattern is:

- Rewrite the self-reference using Longest:  
 $A[i] \leftarrow X \ (A[i] \mid A[i+1]) \ Y$
- Add whole-rule failover to the next-highest relative precedence level using First:  
 $A[i] \leftarrow (X \ (A[i] \mid A[i+1]) \ Y) / A[i+1]$

The highest precedence level  $i = (N-1)$  is assumed to contain a precedence-breaking pattern such as parentheses, therefore should contain a single self-referencing subclause. The rewriting pattern does not need to add a failover clause that matches the next-highest relative precedence level (in fact such a clause would fail to match), and the rewritten

subclause already defers to next highest precedence level (the lowest precedence level  $A[0]$  of the next loop around the precedence hierarchy). Given a highest-precedence rule  $A[N-1] \leftarrow '(' A ')'$ , the rewriting pattern is simply:

- Rewrite the self-reference using Longest:  
 $A[N-1] \leftarrow '(' (A[N-1] \mid A[0]) ')'$

With the grammar rules rewritten into this form, bottom-up parsing no longer gets stuck partway through parsing, while being able to directly build a left-associative structure in the parse tree, including handling multiple loops around the grammar.

The final form of these rules *is still directly left recursive*, but the rule may now be handled correctly by bottom-up parsing. This shows that it was actually not the left recursion that was tripping up bottom-up parsing per se, it was the greediness of the First operator, and the inversion of greediness semantics when switching from top-down to bottom-up.

The reference pika parser (Section 5) detects left-recursive rules, automatically rewriting them into bottom-up-parseable form using this pattern. For example, given the grammar in Listing 16, the reference parser automatically rewrites this grammar into the bottom-up parseable form shown in Listing 17.

---

```

1 E[4] <- '(' E ')';
2 E[3] <- [0-9]+ / [a-z]+;
3 E[2] <- '-' E;
4 E[1,L] <- E op:('*' / '/') E;
5 E[0,L] <- E op:('+ ' / '-') E;
```

---

Listing 16. Example grammar in reference parser notation

---

```

1 E[4] <- '(' (E[4] | E[0]) ')';
2 E[3] <- ([0-9]+ / [a-z]+) / E[4];
3 E[2] <- ('-' (E[2] | E[3])) / E[3];
4 E[1] <- E[1] <- ((E[1] ('*' / '/') E[2]) | (E[2] ('*' / '/') E[2])) / E[2];
5 E[0] <- E[0] <- ((E[0] ('+' / '-') E[1]) | (E[1] ('+' / '-') E[1])) / E[1];
```

---

Listing 17. The result of automatically rewriting the grammar in Listing 16 to work with bottom-up parsing

### 3.5 Automatic conversion of parse tree into an Abstract Syntax Tree (AST)

The structure of the parse tree resulting from a parse is directly induced by the structure of the grammar, i.e. there is one node in the parse tree for each clause and subclause of the grammar that matched the input. Many of these nodes are not needed (for example nodes that match semantically-irrelevant whitespace or comments in the input), and could be suppressed in the output tree. The reference parser allows the user to mark any clause in the grammar with a label, using the syntax (LabelName ':' Clause). After parsing is complete, all matched clauses in the parse tree that do not have an AST node label are simply dropped. The resulting tree is the Abstract Syntax Tree (AST), containing only nodes of interest (Fig. 2).

## 4 BROADER APPLICATION OF INVERSION OF RECURSION

The techniques presented in this paper for converting a top-down recursive algorithm into a bottom-up dynamic programming algorithm may have applications beyond parsing, and could bring the benefits of pika parsing relative to top-down packrat parsing (Section 2.1) to other non-parsing recursive algorithms.



Fig. 2. Automatic conversion of a parse tree into an Abstract Syntax Tree (AST) via AST node labels on grammar clauses

## 5 REFERENCE IMPLEMENTATION

A MIT-licensed reference implementation of the pika parser is available at: <http://github.com/lukehutch/pikaparser>

There are numerous details that can be tricky to get right when implementing this parser, so directly porting the reference parser to other languages will save time compared to writing a parser from scratch using the description in this paper.

The reference implementation includes a *meta-grammar* that is able to parse a PEG grammar in ASCII form, making it easy to write new grammars to be parsed by the reference parser.

## 6 CONCLUSION

The pika parser is a new type of PEG parser that employs dynamic programming to parse a document bottom-up, from characters to the root of the parse tree. Bottom-up parsing enables left-recursive grammars to be parsed directly, making grammar writing simpler, and also enables almost perfect error recovery after syntax errors are encountered, making pika parsers useful for implementing IDEs. The pika parser operates within a moderately small constant performance factor of packrat parsing (both are linear in the length of the input and the depth of the grammar). Additionally, a pika parser can be parallelized to accelerate parsing across available cores, which can compensate for some of the overhead of bottom-up parsing. A lexing step can be applied to reduce the size of the memo table and the amount of work wasted by bottom-up parsing compared to top-down parsing. Mechanisms for implementing precedence and associativity in PEG grammars were demonstrated, and several new insights were provided into the problem of left recursion and the handling of precedence order by PEG parsers, including a fundamental new insight relating the length of a match to its absolute precedence. The mechanism for inverting a top-down dynamic recursive algorithm to produce a bottom-up dynamic programming algorithm may have broader applications. A reference implementation of the pika parser is available under the MIT license.

## REFERENCES

- [1] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. 2006. *Compilers: principles, techniques, and tools*.
- [2] Alfred V Aho and Jeffrey D Ullman. 1972. *The theory of parsing, translation, and compiling*. Vol. 1. Prentice-Hall Englewood Cliffs, NJ.
- [3] Alfred V. Aho and Jeffrey D. Ullman. 1977. *Principles of Compiler Design*. Addison-Wesley.
- [4] Noam Chomsky. 1956. Three models for the description of language. *IRE Transactions on information theory* 2, 3 (1956), 113–124.
- [5] Noam Chomsky. 1959. On certain formal properties of grammars. *Information and control* 2, 2 (1959), 137–167.
- [6] William John Cocke, Jacob T Schwartz, and Courant Institute of Mathematical Sciences. 1970. *Programming languages and their compilers*. Courant Institute of Mathematical Sciences.
- [7] H Daniel. 1967. Younger. Recognition and parsing of context-free languages in time  $n^3$ . *Information and control* 10, 2 (1967), 189–208.
- [8] Sérgio Queiroz de Medeiros, Gilney de Azevedo Alvez Junior, and Fabio Mascarenhas. 2020. Automatic syntax error reporting and recovery in parsing expression grammars. *Science of Computer Programming* 187 (2020), 102373.
- [9] Sérgio Queiroz de Medeiros and Fabio Mascarenhas. 2018. Syntax error recovery in parsing expression grammars. In *Proceedings of the 33rd Annual ACM Symposium on Applied Computing*. 1195–1202.
- [10] Franklin Lewis DeRemer. 1969. *Practical translators for LR (k) languages*. Ph.D. Dissertation. Massachusetts Institute of Technology.
- [11] Bryan Ford. 2002. *Packrat parsing: a practical linear-time algorithm with backtracking*. Ph.D. Dissertation. Massachusetts Institute of Technology.
- [12] Bryan Ford. 2004. Parsing expression grammars: a recognition-based syntactic foundation. In *Proceedings of the 31st ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. 111–122.
- [13] Richard A Frost and Rahmatullah Hafiz. 2006. A new top-down parsing algorithm to accommodate ambiguity and left recursion in polynomial time. *ACM SIGPLAN Notices* 41, 5 (2006), 46–54.
- [14] Richard A. Frost, Rahmatullah Hafiz, and Paul C. Callaghan. 2007. Modular and Efficient Top-down Parsing for Ambiguous Left-Recursive Grammars. In *Proceedings of the 10th International Conference on Parsing Technologies (Prague, Czech Republic) (IWPT '07)*. Association for Computational Linguistics, USA, 109–120.
- [15] Tadao Kasami. 1966. An efficient recognition and syntax-analysis algorithm for context-free languages. *Coordinated Science Laboratory Report no. R-257* (1966).
- [16] Stephen Cole Kleene. 1951. Representation of events in nerve nets and finite automata. *RAND Research Memorandum RM-704* (1951).
- [17] Donald E Knuth. 1962. History of writing compilers. In *Proceedings of the 1962 ACM National Conference on Digest of Technical Papers*. 43.
- [18] Donald E Knuth. 1965. On the translation of languages from left to right. *Information and control* 8, 6 (1965), 607–639.
- [19] Joseph B Kruskal. 1983. An overview of sequence comparison: Time warps, string edits, and macromolecules. *SIAM review* 25, 2 (1983), 201–237.
- [20] Bernard Lang. 1974. Deterministic techniques for efficient non-deterministic parsers. In *International Colloquium on Automata, Languages, and Programming*. Springer, 255–269.

- [21] Robert McNaughton and Hisao Yamada. 1960. Regular expressions and state graphs for automata. *IRE transactions on Electronic Computers* 1 (1960), 39–47.
- [22] Peter Norvig. 1991. Techniques for automatic memoization with applications to context-free parsing. *Computational Linguistics* 17, 1 (1991), 91–98.
- [23] Klaus Samelson and Friedrich L Bauer. 1960. Sequential formula translation. *Commun. ACM* 3, 2 (1960), 76–83.
- [24] Michael Sipser. 2012. *Introduction to the Theory of Computation*. Cengage Learning.
- [25] Masaru Tomita. 2013. *Efficient parsing for natural language: a fast algorithm for practical systems*. Vol. 8. Springer Science & Business Media.
- [26] Alessandro Warth, James R Douglass, and Todd Millstein. 2008. Packrat parsers can support left recursion. In *Proceedings of the 2008 ACM SIGPLAN symposium on Partial evaluation and semantics-based program manipulation*. 103–110.