

Parsing Expression Grammar and Packrat Parsing (Survey)

IPLAS Seminar Oct 27, 2009

Kazuhiro Inaba

This Talk is Based on These Resources

- ▶ The Packrat Parsing and PEG Page (by Bryan Ford)
 - ▶ <http://pdos.csail.mit.edu/~baford/packrat/>
 - ▶ (was active till early 2008)
 - ▶ A. Birman & J. D. Ullman, "Parsing Algorithms with Backtrack", Information and Control (23), 1973
 - ▶ B. Ford, "Packrat Parsing: Simple, Powerful, Lazy, Linear Time", ICFP 2002
 - ▶ B. Ford, "Parsing Expression Grammars: A Recognition-Based Syntactic Foundation", POPL 2004
-

Outline

- ▶ **What is PEG?**
 - ▶ Introduce the core idea of Parsing Expression Grammars
 - ▶ **Packrat Parsing**
 - ▶ Parsing Algorithm for the core PEG
 - ▶ **Packrat Parsing Can Support More...**
 - ▶ Syntactic predicates
 - ▶ **Full PEG**
 - ▶ This is what is called “PEG” in the literature.
 - ▶ **Theoretical Properties of PEG**
 - ▶ **PEG in Practice**
-

What is PEG?

- ▶ Yet Another Grammar Formalism
 - ▶ Intended for describing grammars of programming languages (not for NL, nor for program analysis)
 - ▶ As simple as Context-Free Grammars
 - ▶ Linear-time parsable
 - ▶ Can express:
 - ▶ All deterministic CFLs (LR(k) languages)
 - ▶ Some non-CFLs
-

What is PEG? – Comparison to CFG

(Predicate-Free) Parsing Expression Grammar

- ▶ $A \leftarrow B C$
 - ▶ Concatenation
- ▶ $A \leftarrow B / C$
 - ▶ Prioritized Choice
 - ▶ When both B and C matches, prefer B

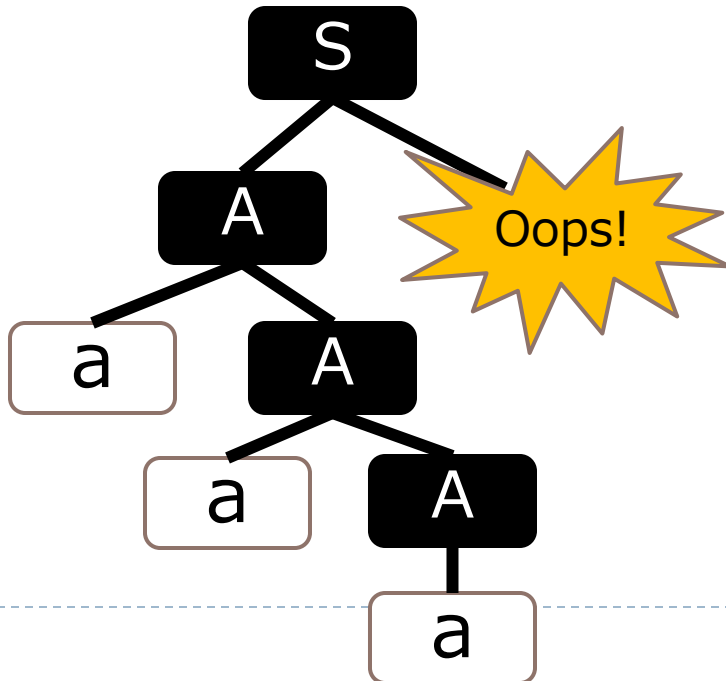
Context-Free Grammar

- ▶ $A \rightarrow B C$
 - ▶ Concatenation
 - ▶ $A \rightarrow B \mid C$
 - ▶ Unordered Choice
 - ▶ When both B and C matches, either will do
-

Example

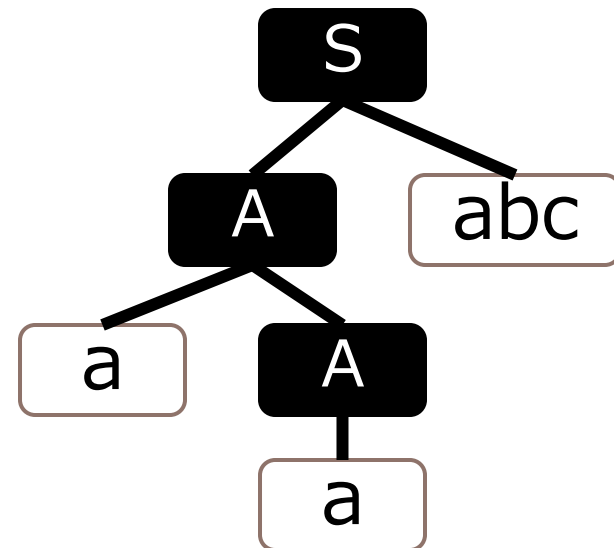
(Predicate-Free) Parsing Expression Grammar

- ▶ $S \leftarrow A a b c$
- ▶ $A \leftarrow a A / a$
 - ▶ S fails on "aaabc".



Context-Free Grammar

- ▶ $S \rightarrow A a b c$
- ▶ $A \rightarrow a A \mid a$
 - ▶ S recognizes "aaabc".



Another Example

(Predicate-Free) Parsing Expression Grammar

- ▶ $S \leftarrow E ;$
 - / while (E) S
 - / if (E) S else S
 - / if (E) S
 - / ...
- ▶

```
if(x>0)
  if(x<9)
    y=1;
  else
    y=3; unambiguous
```

Context-Free Grammar

- ▶ $S \rightarrow E ;$
 - | while (E) S
 - | if (E) S else S
 - | if (E) S
 - | ...
- ▶

```
if(x>0)
  if(x<9)
    y=1;
  else
    y=3;   ambiguous
```

Formal Definition

- ▶ Predicate-Free PEG G is $\langle N, \Sigma, S, R \rangle$
 - ▶ N : Finite Set of Nonterminal Symbols
 - ▶ Σ : Finite Set of Terminal Symbols
 - ▶ $S \in N$: Start Symbol
 - ▶ $R \in N \rightarrow \text{rhs}$: Rules, where
 - ▶ $\text{rhs} ::= \varepsilon$
 - ▶ $A \quad (A \in N)$
 - ▶ $a \quad (a \in \Sigma)$
 - ▶ rhs / rhs
 - ▶ $\text{rhs} \text{ rhs}$
 - ▶ Note: $A \leftarrow \text{rhs}$ stands for $R(A) = \text{rhs}$
 - ▶ Note: Left-recursion is not allowed
-

Semantics

- ▶ $[[e]] :: \text{String} \rightarrow \text{Maybe String}$ where $\text{String} = \Sigma^*$
 - ▶ $[[c]] = \lambda s \rightarrow \text{case } s \text{ of}$ (for $c \in \Sigma$)
 - ▶ $c : t \rightarrow \text{Just } t$
 - ▶ $_ \rightarrow \text{Nothing}$
 - ▶ $[[e1 e2]] = \lambda s \rightarrow \text{case } [[e1]] s \text{ of}$
 - ▶ $\text{Just } t \rightarrow [[e2]] t$
 - ▶ $\text{Nothing} \rightarrow \text{Nothing}$
 - ▶ $[[e1 / e2]] = \lambda s \rightarrow \text{case } [[e1]] s \text{ of}$
 - ▶ $\text{Just } t \rightarrow \text{Just } t$
 - ▶ $\text{Nothing} \rightarrow [[e2]] s$
 - ▶ $[[\varepsilon]] = \lambda s \rightarrow \text{Just } s$
 - ▶ $[[A]] = [[R(A)]]$ (recall: $R(A)$ is the unique rhs of A)
-

Example (Complete Consumption)

$S \leftarrow a S b / c$

- ▶ $[[S]]$ "acb" = Just ""
 - ▶ $[[aSb]]$ "acb" = Just ""
 - ▶ $[[a]]$ "acb" = Just "cb"
 - ▶ $[[S]]$ "cb" = Just "b"
 - $[[aSb]]$ "cb" = Nothing
 - ▶ $[[a]]$ "cb" = Nothing
 - $[[c]]$ "cb" = Just "b"
 - ▶ $[[b]]$ "b" = Just ""
-

Example (Failure, Partial Consumption)

$S \leftarrow a S b / c$

- ▶ $[[S]] \text{ "b"}$ = Nothing
- ▶ $[[aSb]] \text{ "b"}$ = Nothing
 - ▶ $[[a]] \text{ "b"}$ = Nothing
- ▶ $[[c]] \text{ "b"}$ = Nothing

-
- ▶ $[[S]] \text{ "cb"}$ = Just "b"
 - ▶ $[[aSb]] \text{ "cb"}$ = Nothing
 - ▶ $[[a]] \text{ "cb"}$ = Nothing
 - ▶ $[[c]] \text{ "cb"}$ = Just "b"

Example (Prioritized Choice)

$S \leftarrow A a$

$A \leftarrow a A / a$

- ▶ $[[S]]$ "aa" = Nothing
 - ▶ Because $[[A]]$ "aa" = Just "", not Just "a"
 - ▶ $[[A]]$ "aa" = Just ""
 - ▶ $[[a]]$ "aa" = Just "a"
 - ▶ $[[A]]$ "a" = Just ""
 - ...
-

“Recognition-Based”

- ▶ In “generative” grammars such as CFG, each nonterminal defines a language (set of strings) that it generates.
 - ▶ In “recognition-based” grammars, each nonterminal defines a parser (function from string to something) that it recognizes.
-

Outline

- ▶ What is PEG?
 - ▶ Introduce the core idea of Parsing Expression Grammars
 - ▶ Packrat Parsing
 - ▶ Parsing Algorithm for the core PEG
 - ▶ Packrat Parsing Can Support More...
 - ▶ Syntactic predicates
 - ▶ Full PEG
 - ▶ This is what is called “PEG” in the literature.
 - ▶ Theoretical Properties of PEG
 - ▶ PEG in Practice
-

Parsing Algorithm for PEG

- ▶ Theorem: Predicate-Free PEG can be parsed in **linear time** wrt the length of the input string.
 - ▶ Proof
 - ▶ By Memoization
 - (All arguments and outputs of
 - $[[e]] :: \text{String} \rightarrow \text{Maybe String}$are the suffixes of the input string)
-

Parsing Algorithm for PEG

- ▶ How to Memoize?
 - ▶ Tabular Parsing [Birman&Ullman73]
 - ▶ Prepare a table of size $|G| \times |\text{input}|$, and fill it from right to left.
 - ▶ Packrat Parsing [Ford02]
 - ▶ Use lazy evaluation.
-

Parsing PEG (1: Vanilla Semantics)

$S \leftarrow aS / a$

- ▶ `doParse = parseS :: String -> Maybe String`
 - ▶ `parseA s =`
 - ▶ `case s of 'a':t -> Just t`
 - ▶ `_ -> Nothing`
 - ▶ `parseS s = alt1 `mplus` alt2 where`
 - ▶ `alt1 = case parseA s of`
 - ▶ `Just t -> case parseS t of`
 - `Just u -> Just u`
 - `Nothing -> Nothing`
 - ▶ `Nothing -> Nothing`
 - ▶ `alt2 = parseA s`
-

Parsing PEG (2: Valued)

$S \leftarrow aS / a$

- ▶ `doParse = parseS :: String -> Maybe (Int, String)`
- ▶ `parseA s =`
 - ▶ `case s of 'a':t -> Just (1, t)`
 - ▶ `_ -> Nothing`
- ▶ `parseS s = alt1 `mplus` alt2` where
 - ▶ `alt1 = case parseA s of`
 - ▶ `Just (n,t) -> case parseS t of`
 - `Just (m,u) -> Just (n+m,u)`
 - `Nothing -> Nothing`
 - ▶ `Nothing -> Nothing`
 - ▶ `alt2 = parseA s`

Parsing PEG (3: Packrat Parsing)

$S \leftarrow aS / a$

- ▶ type Result = Maybe (Int, Deriv)
- ▶ data Deriv = D Result Result

- ▶ doParse :: String -> Deriv
- ▶ doParse s = d where
 - ▶ d = D resultS resultA
 - ▶ resultS = parseS d
 - ▶ resultA = case s of 'a':t -> Just (1,next)
 - ▶ _ -> Nothing
 - ▶ next = doParse (tail s)

▶ ...

Parsing PEG (3: Packrat Parsing, cnt'd)

$S \leftarrow aS / a$

- ▶ type Result = Maybe (Int, Deriv)
- ▶ data Deriv = D Result Result
- ▶ parseS :: Deriv -> Result
- ▶ parseS (D rS0 rA0) = alt1 `mplus` alt2 where
 - ▶ alt1 = case rA0 of
 - ▶ Just (n, D rS1 rA1) -> case rS1 of
 - Just (m, d) -> Just (n+m, d)
 - Nothing -> Nothing
 - ▶ Nothing -> Nothing
 - ▶ alt2 = rA0

- ▶ alt1 = case parseA s of
 - ▶ Just (n,t)-> case parseS t of
 - Just (m,u)-> Just (n+m,u)
 - Nothing -> Nothing
 - ▶ Nothing -> Nothing
- ▶ alt2 = parseA s

Packrat Parsing Can Do More

- ▶ Without sacrificing linear parsing-time, more operators can be added. Especially, “syntactic predicates”:
 - ▶ $[[\&e]] = \lambda s \rightarrow \text{case } [[e]] \text{ s of}$
 - ▶ Just $_ \rightarrow \text{Just } s$
 - ▶ Nothing $\rightarrow \text{Nothing}$
 - ▶ $[[!e]] = \lambda s \rightarrow \text{case } [[e]] \text{ s of}$
 - ▶ Just $_ \rightarrow \text{Nothing}$
 - ▶ Nothing $\rightarrow \text{Just } s$
-

Formal Definition of PEG

▶ PEG G is $\langle N, \Sigma, S, R \in N \rightarrow \text{rhs} \rangle$ where

- ▶ $\text{rhs} ::= \varepsilon$
 - | $A \quad (\in N)$
 - | $a \quad (\in \Sigma)$
 - | rhs / rhs
 - | $\text{rhs} \text{ rhs}$
 - | $\&\text{rhs}$
 - | $!\text{rhs}$
 - | $\text{rhs}?$ (eqv. to X where $X \leftarrow \text{rhs} / \varepsilon$)
 - | rhs^* (eqv. to X where $X \leftarrow \text{rhs} X / \varepsilon$)
 - | rhs^+ (eqv. to X where $X \leftarrow \text{rhs} X / \text{rhs}$)
-

Example: A Non Context-Free Language

▶ $\{a^n b^n c^n \mid n > 0\}$

is recognized by

- ▶ $S \leftarrow \&X a^* Y !a !b !c$
 - ▶ $X \leftarrow aXb / ab$
 - ▶ $Y \leftarrow bYc / bc$
-

Example: C-Style Comment

- ▶ C-Style Comment
 - ▶ **Comment** ← `/* ((! */) Any)* */`
 - ▶ (for readability, meta-symbols are **colored**)
 - ▶ Though this is a regular language, it cannot be written this easy in conventional regex.
-

Outline

- ▶ **What is PEG?**
 - ▶ Introduce the core idea of Parsing Expression Grammars
 - ▶ **Packrat Parsing**
 - ▶ Parsing Algorithm for the core PEG
 - ▶ **Packrat Parsing Can Support More...**
 - ▶ Syntactic predicates
 - ▶ **Full PEG**
 - ▶ This is what is called “PEG” in the literature.
 - ▶ **Theoretical Properties of PEG**
 - ▶ **PEG in Practice**
-

Theoretical Properties of PEG

- ▶ Two Topics

- ▶ Properties of Languages Defined by PEG
- ▶ Relationship between PEG and predicate-free PEG

Language Defined by PEG

- ▶ For a parsing expression e
 - ▶ [Ford04] $F(e) = \{w \in \Sigma^* \mid [[e]]w \neq \text{Nothing}\}$
 - ▶ [BU73] $B(e) = \{w \in \Sigma^* \mid [[e]]w = \text{Just ""}\}$
 - ▶ [Redziejowski08]
 - ▶ R. R. Redziejowski, "Some Aspects of Parsing Expression Grammar", *Fundamenta Informaticae*(85), 2008
 - ▶ Investigation on concatenation $[[e_1 e_2]]$ of two PEGs
 - ▶ $S(e) = \{w \in \Sigma^* \mid \exists u. [[e]]wu = \text{Just } u\}$
 - ▶ $L(e) = \{w \in \Sigma^* \mid \forall u. [[e]]wu = \text{Just } u\}$
-

Properties of $F(e) = \{w \in \Sigma^* \mid [[e]]w \neq \text{Nothing}\}$

- ▶ $F(e)$ is context-sensitive
 - ▶ Contains all deterministic CFL
 - ▶ Trivially Closed under Boolean Operations
 - ▶ $F(e_1) \cap F(e_2) = F(\&e_1 e_2)$
 - ▶ $F(e_1) \cup F(e_2) = F(e_1 / e_2)$
 - ▶ $\sim F(e) = F(!e)$
 - ▶ Undecidable Problems
 - ▶ " $F(e) = \Phi$ "? is undecidable
 - ▶ Proof is similar to that of intersection emptiness of context-free languages
 - ▶ " $F(e) = \Sigma^*$ "? is undecidable
 - ▶ " $F(e_1) = F(e_2)$ "? is undecidable
-

Properties of $B(e) = \{w \in \Sigma^* \mid [[e]]w = \text{Just } \text{""}\}$

- ▶ $B(e)$ is context-sensitive
 - ▶ Contains all deterministic CFL
 - ▶ For predicate-free e_1, e_2
 - ▶ $B(e_1) \cap B(e_2) = B(e_3)$ for some predicate-free e_3
 - ▶ For predicate-free & well-formed e_1, e_2 where well-formed means that $[[e]] s$ is either $\text{Just } \text{""}$ or Nothing
 - ▶ $B(e_1) \cup B(e_2) = B(e_3)$ for some pf&wf e_3
 - ▶ $\sim B(e_1) = B(e_3)$ for some predicate-free e_3
 - ▶ Emptiness, Universality, and Equivalence is undecidable
-

Properties of $B(e) = \{w \in \Sigma^* \mid [[e]]w = \text{Just } \text{“”}\}$

- ▶ Forms AFDL, i.e.,
 - ▶ $\text{markedUnion}(L_1, L_2) = aL_1 \cup bL_2$
 - ▶ $\text{markedRep}(L_1) = (aL_1)^*$
 - ▶ **marked inverse GSM** (inverse image of a string transducer with explicit endmarker)
 - ▶ [Chandler69] AFDL is closed under many other operations, such as **left-/right- quotients**, **intersection with regular sets**, ...
 - ▶ W. J. Chandler, “Abstract Families of Deterministic Languages”, STOC 1969
-

Predicate Elimination

▶ Theorem: $G = \langle N, \Sigma, S, R \rangle$ be a PEG such that $F(S)$ does not contain ε . Then there is an equivalent predicate-free PEG.

▶ Proof (Key Ideas):

▶ $[[\&e]] = [[!!e]]$

▶ $[[!e C]] = [[(e Z / \varepsilon) C]]$ for ε -free C

▶ where $Z = (\sigma_1 / \cdots / \sigma_n) Z / \varepsilon$, $\{\sigma_1, \dots, \sigma_n\} = \Sigma$

Predicate Elimination

▶ Theorem: PEG is strictly more powerful than predicate-free PEG

▶ Proof:

▶ We can show, for predicate-free e ,

▶ $\forall w. ([[e]] \text{ ""} = \text{Just ""} \iff [[e]] w = \text{Just } w)$

by induction on $|w|$ and on the length of derivation

▶ Thus we have

▶ $\text{""} \in F(S) \iff F(S) = \Sigma^*$

but this is not the case for general PEG (e.g., $S \leftarrow !a$)

Outline

- ▶ **What is PEG?**
 - ▶ Introduce the core idea of Parsing Expression Grammars
 - ▶ **Packrat Parsing**
 - ▶ Parsing Algorithm for the core PEG
 - ▶ **Packrat Parsing Can Support More...**
 - ▶ Syntactic predicates
 - ▶ **Full PEG**
 - ▶ This is what is called “PEG” in the literature.
 - ▶ **Theoretical Properties of PEG**
 - ▶ **PEG in Practice**
-

PEG in Practice

- ▶ Two Topics
 - ▶ When is PEG useful?
 - ▶ Implementations
-

When is PEG useful?

- ▶ When you want to unify lexer and parser
 - ▶ For packrat parsers, it is easy.
 - ▶ For LL(1) or LALR(1) parsers, it is not.

```
list<list<string>>>
```

- ▶ Error in C++98, because >> is RSHIFT, not two closing angle brackets
- ▶ Ok in Java5 and C++1x, but with strange grammar

```
(* nested (* comment *) *)
```

```
s = "embedded code #{1+2+3} in string"
```

Implementations

- **Java:**
 - [Rats!](#) by [Robert Grimm](#), a part of the [eXI](#)
 - [ANTLR](#), a well-established parser generator packrat parsing with LL parsing technique
 - [LGI](#), a dynamic PEG-based parser generator
 - **Python:**
 - The [pyparsing](#) monadic parsing combinator
 - [Packrat parsing support](#) has also been included
 - **Haskell:**
 - [Frisby](#) by [John Meacham](#) is a monadic parser support dynamic specification of composition
 - [Pappy](#) by [Bryan Ford](#) is a simple prototype
 - **C, C++:**
 - The [Narwhal](#) compiler suite by Gordon Tis
 - The [PEG Template Library](#) for [C++0x](#) by [T](#)
 - The [peg/leg](#) parser generator emphasizes
 - **C#:** [NPEG](#) is a library providing objects to build
 - **JavaScript:** [OMeta](#) supports PEG-based patterns
 - **Tcl:** The new [grammar::peg](#) module supports grammar
 - **Smalltalk:** [OMeta](#) provides PEG-based patterns
 - **Scheme:** Tony Garnock-Jones has written [a parser](#)
 - **Common Lisp:** [CL-peg](#) by John Leuner supports
 - **Lua:** Roberto Ierusalimschy has provided [the LI](#)
 - **Ruby** now has the [Treetop](#) grammar description
-

Performance (Rats!)

- ▶ R. Grimm, “Better Extensibility through Modular Syntax”, PLDI 2006
 - ▶ Parser Generator for PEG, used, e.g., for Fortress

System	Algorithm	Modules	Lex	AST	LoC
<i>Rats!</i>	PEG	9	—	—	790
SDF2	GLR	57	—	—	1,680
Elkhound	LALR/GLR	1	1	1	2,370
ANTLR	LL	1	1	—	1,280
JavaCC	LL	1	1	—	1,240

System	Recognizer		Parser	
	T-put	Heap Util.	T-put	Heap Util.
<i>Rats!</i>	518.0	51.5	317.0	58.0
SDF2	136.1	—	21.4	—
Elkhound	141.5	—	139.4	—
ANTLR	538.6	11.5	393.6	28.0
JavaCC	1,114.3	10.6	382.9	63.2

Experiments
on Java1.4
grammar,
with sources
of size
0.7 ~ 70KB

PEG in Fortress Compiler

- ▶ Syntactic Predicates are widely used
 - ▶ (though I'm not sure whether it is essential, due to my lack of knowledge on Fortress...)

```
/* The operator "|->" should not be in the left-hand sides of map
expressions and map/array comprehensions.
```

```
*/
```

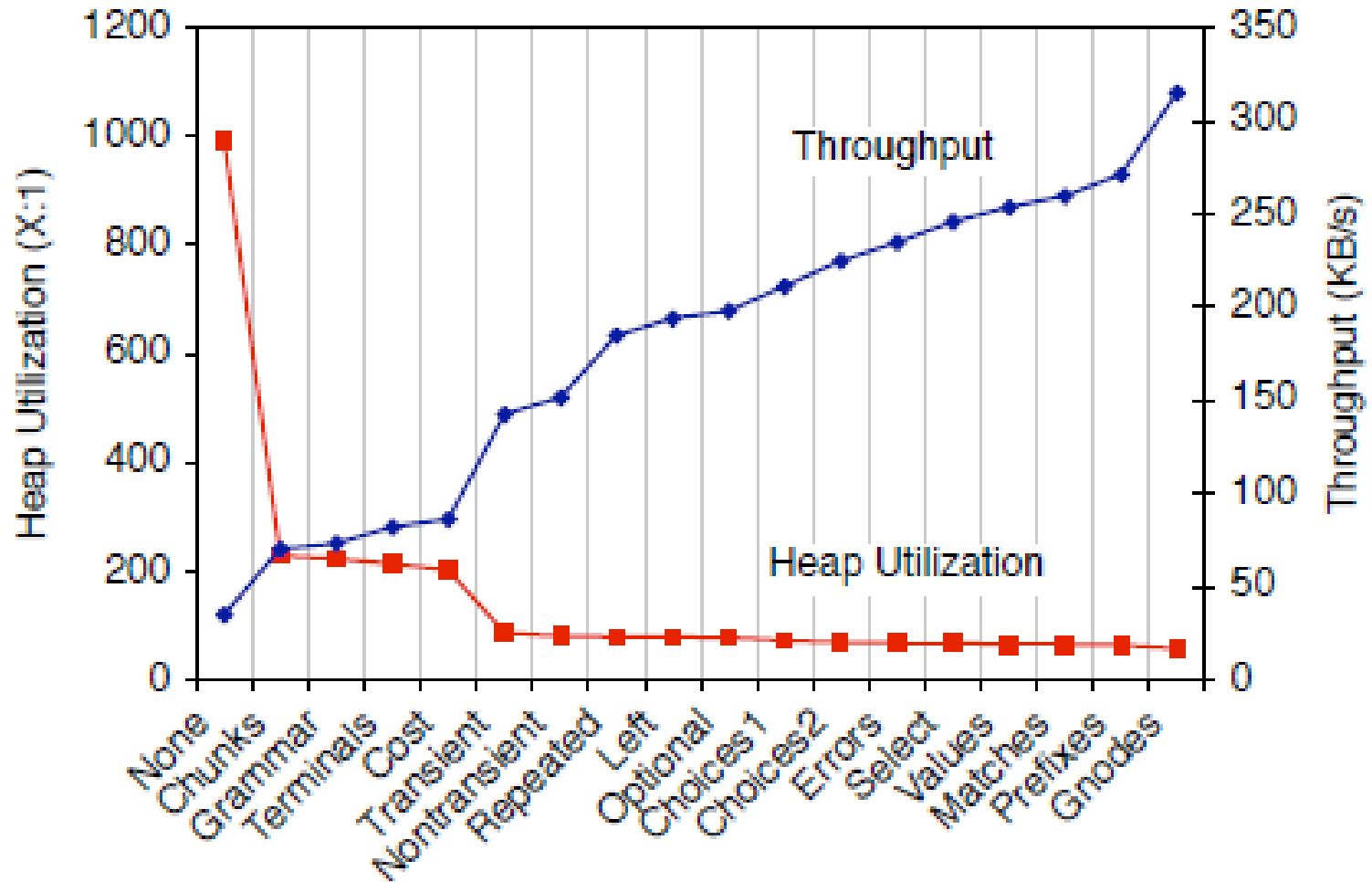
```
String mapstoOp =
```

```
!("|->" w Expr (w mapsto / wr bar / w closecurly / w comma)) "|->" ;
```

```
/* The operator "<-" should not be in the left-hand sides of
generator clause lists. */
```

```
String leftarrowOp = !("<-" w Expr (w leftarrow / w comma)) "<-";
```

Optimizations in Rats!



Summary

- ▶ Parsing Expression Grammar (PEG) ...
 - ▶ has prioritized choice e_1/e_2 , rather than unordered choice $e_1|e_2$.
 - ▶ has syntactic predicates $\&e$ and $!e$, which can be eliminated if we assume ε -freeness.
 - ▶ might be useful for unified lexer-parser.
 - ▶ can be parsed in $O(n)$ time, by memoizing.
-