

# PEGs and Packrat Parsing

April 13, 2011

CS 132: Compiler Design

## CHOMSKY HIERARCHY

- ✓ Type 0: Unrestricted Grammars
- ✓ Type 1: Context-Sensitive Grammars
- ✓ Type 2: Context-Free Grammars
- ✓ Type 3: Regular Grammars

The formalism seems backwards if we want to recognize (parse) strings in the language...

# PARSING EXPRESSION GRAMMARS (PEGs)

*Recognition-Based, rather than generative.*

Define nonterminals as in CFGs, but right-hand-side can involve:

$e?$	Optional (greedy)
$e^*$	Zero-or-more (greedy)
$e^+$	One-or-more (greedy)
$\&e$	Matches (does not consume)
$!e$	Doesn't Match (does not consume)
$e_1 e_2$	Sequencing
$e_1 / e_2$	Choice (prioritized)

Advantages:

- ✓ Deterministic by definition
- ✓ Linear-time parsing

## SIMPLE PEG EXAMPLES

A  $\leftarrow$  a b / a

B  $\leftarrow$  a / a b

C  $\leftarrow$  a a\*

D  $\leftarrow$  a\* a

e?	Optional ( <i>greedy</i> )
e*	Zero-or-more ( <i>greedy</i> )
e <sup>+</sup>	One-or-more ( <i>greedy</i> )
&e	Matches ( <i>does not consume</i> )
!e	Doesn't Match ( <i>does not consume</i> )
e <sub>1</sub> e <sub>2</sub>	Sequencing
e <sub>1</sub> /e <sub>2</sub>	Choice ( <i>prioritized</i> )

## RESTRICTION: LEFT RECURSION

Normally, PEGs cannot be left-recursive:

$$E \leftarrow E + E / n$$

$$F \leftarrow n / E + E$$

(Recent work lifts this restriction, but we lose the linear-time-parsing guarantee.)

## UNIFIED LEXING AND PARSING

PEGs commonly combine lexing and parsing in a single grammar.  
(Regex operators make this less painful.)

Advantages:

- ✓ Everything in one place with one formalism
- ✓ No problem with context-dependent lexing, as in
  - ▶ `vector<vector<int>> v(n>>2);`

Disadvantages:

- ✓ Need to specify whitespace in the PEG  
(E.g., after every lexical token)

Identifier ← Letter (Letter / Number)\* Whitespace

Letter ← [A-Za-z]

Number ← [0-9]

Whitespace ← ([ \t\n] / Comment)\*

## LEXICAL EXAMPLES

### Character literals

```
Literal ← ['] (!['] Char)* [']  
Char    ← [\\] [nrt'"\\[\\]\\]  
        / ![\\] .
```

### Haskell comments

```
Comment ← "{-" (Comment / !"-}" .)* "-}"
```

### End-of-file

```
EndOfFile ← !.
```

## PARSING EXAMPLES

(Whitespace OMITTED)

Expressions:

$E \leftarrow T ( "+" E / "-" E )^*$

$T \leftarrow F ( "*" F / "/" F )^*$

$F \leftarrow \text{Number} / "(" E ")"$

Unambiguous if-else:

$\text{Statement} \leftarrow \text{"if" "(" Expression ")" Statement "else" Statement}$   
 $\quad / \text{"if" "(" Expression ")" Statement}$

## A NON-CONTEXT-FREE PEG

$A \leftarrow a A b \ / \ \epsilon$

$C \leftarrow b C c \ / \ \epsilon$

$S \leftarrow \&(A \ !b) a^* C \ !.$

## PACKRAT PARSING

Intuition: Follow the grammar (like recursive-descent), but:

1. Memoize: “Have I tried to match this nonterminal at this point in the string before? If so, how much of the input did it match?”
2. Backtracking is limited in PEGs;  
Each nonterminal can match a prefix of a string in exactly one way.



## APPLICATION: FORTRESS PROGRAMMING LANGUAGE

- ✓ Designed for high-performance computing
- ✓ Developed by Guy Steele and others at Sun Labs (now Oracle)
- ✓ Relevance here: mathematical *and extensible* syntax.

```
cgit = 25
z := Vector[\E,n\](0)
r := x
rho := r DOT r
p := r
for j <- seq(1#cgit) do
  q = A p
  alpha = rho/(p DOT q)
  z += alpha p
  rho0 = rho
  r -= alpha q
  rho := r DOT r
  beta = rho / rho0
  p := r + beta p
end
(z, ||x - A z||)
```

```
cgit = 25
z := Vector[[E,n]](0)
r := x
ρ := r · r
p := r
for j ← seq(1 # cgit) do
  q = A p
  α =  $\frac{\rho}{p \cdot q}$ 
  z += α p
  ρ0 = ρ
  r -= α q
  ρ := r · r
  β = ρ/ρ0
  p := r + β p
end
(z, ||x - A z||)
```

## EXTENSIBLE SYNTAX

The Fortress programming language is intended to grow over time to accommodate the changing needs of its users [4]. One mechanism to allow for such growth is syntactic abstraction: It is possible to add new syntactic constructs to the language in libraries, defining new constructs in terms of old ones. In this manner, the language can gracefully adapt to unanticipated needs as they become apparent. Parsing of new constructs can be done alongside parsing of primitive constructs, allowing programmers to detect syntax errors in use sites of new constructs early. Programs in domain-specific languages can be embedded in Fortress programs and parsed along with their host programs. Moreover, the defini-

From Allen et al., "Growing a Syntax," 2009.

## EXTENSIBLE SYNTAX

The parsing stage relies on *Parsing Expression Grammars* (PEGs) [16] because they have some advantages over usual Context-Free Grammars (CFG) [10] based parsing formalisms such as LL and LALR(k). PEGs are unambiguous, are closed under union, and integrates lexing with parsing. We need the closure property because we want to combine PEGs for different grammars and the integrated lexing and parsing makes it easy to achieve our goal of “similar syntax for definition and use”. Furthermore, the parsers based on PEGs allow a linear execution time compared to the generalized CFG parsers. We briefly introduce PEGs in Section 5.1.1 along with a description of our pattern language which is effectively a variant of *Parsing Expressions*.

From Allen et al., “Growing a Syntax,” 2009.

## EXTENSIBLE SYNTAX

We have evaluated the design of our syntactic abstraction system by implementing the `for` loop example shown in Figure 3, and also a grammar that recognizes regular expressions and one recognizing XML. These examples show that the system is suitable for language extensions and domain specific languages by allowing their care-free implementation in general. The grammars for XML and regular expressions can be found as part of the open-source Fortress interpreter [3].

The `for` loop example demonstrates a weakness of the previous

From Allen et al., “Growing a Syntax,” 2009.