

# A Partial Conversion of Parsing Expression Grammars to Deterministic Finite Automata

Nariyoshi Chida and Kimio Kuramitsu  
Yokohama National University, Japan

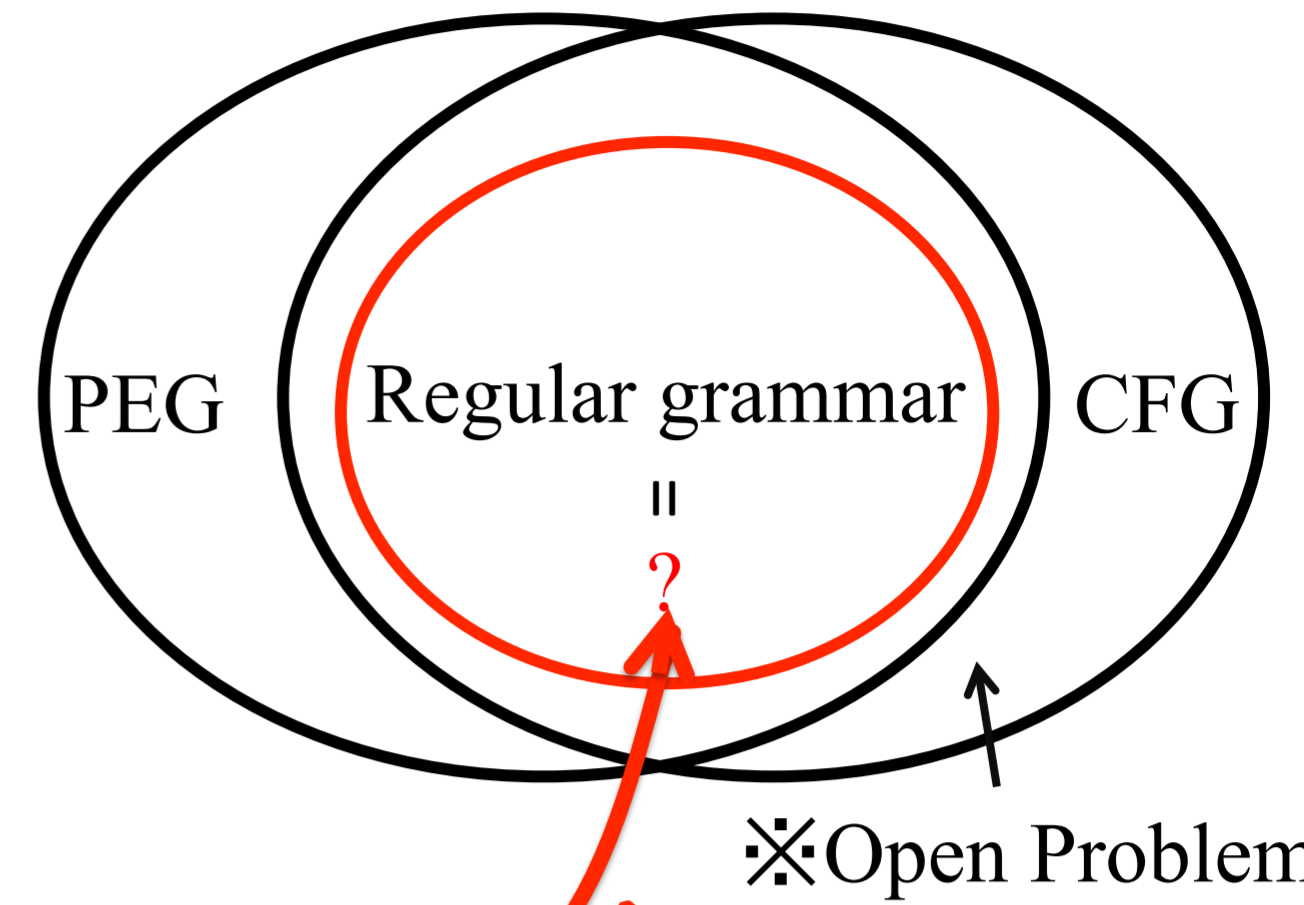


Visit the website! → <http://regex-and-pe-to-dfa.com>

## Background : Parsing Expression Grammars (PEGs)

PEGs are formal grammars introduced by Ford in 2004.

- Features of PEGs
  - Parsed in linear time using a memorizing parser called *packrat parser*
  - Easily implemented by recursive descent parsers
    - There is no need to implement a lexer (i.e. PEGs are scanner-less)
  - Deterministic
    - Behavior of each operators are deterministic (e.g. repetition operators are greedy)
  - Recognize languages that is not context-free such as  $\{a^n b^n c^n | n \geq 0\}$



## What is a subclass of PEGs that is equivalent to DFAs?

Since PEGs are relatively new, there are several unsolved problems.

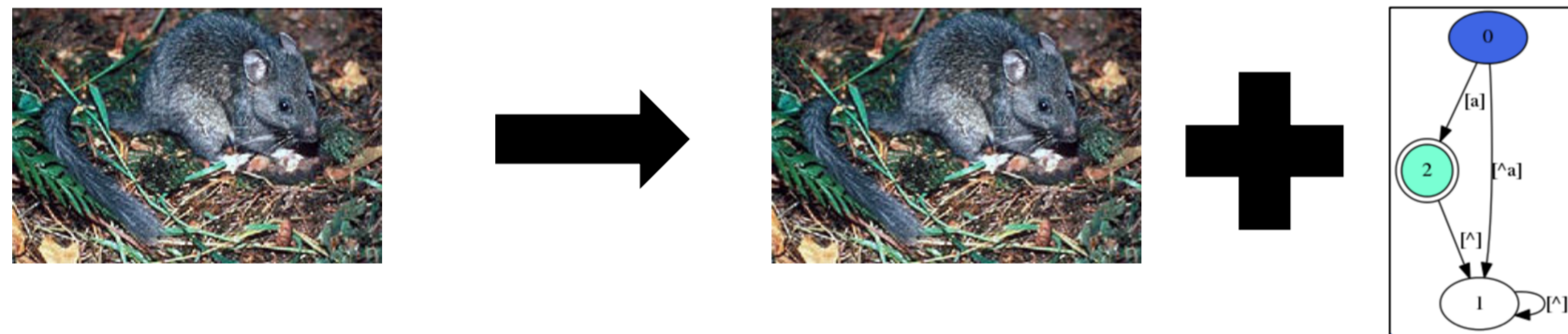
One of the problems is that the revealing a subclass of PEGs that is equivalent to DFAs.

## Motivation :

### Eliminating backtrackings

Backtracking arises when packrat parsers fail to match.

- Regular Grammar in PEGs
  - This allows to apply some techniques from the theory of RG to PEGs.
  - For example, the techniques of the theory of RG are used for faster parsing
    - A DFA optimization for PEG-based parsers **Improved!**



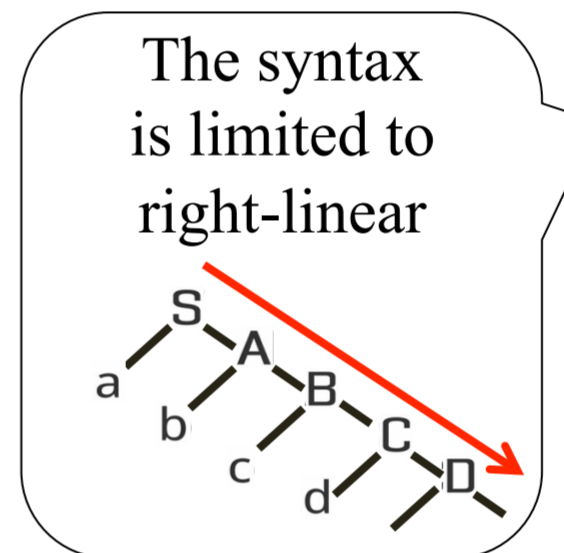
We formalize the subclass as linear PEG (LPEG).

## Idea : Linear PEGs

A LPEG is a set of named *linear parsing expressions*. They are specified by rules of the form  $A \leftarrow e$  where  $e$  is a linear parsing expression and  $A$  is the name. The syntax of  $e$  is shown in Figure 1.

- Example 1 (LPEG)

$A \leftarrow a A / b B / c$   
 $B \leftarrow a B / b A / c$



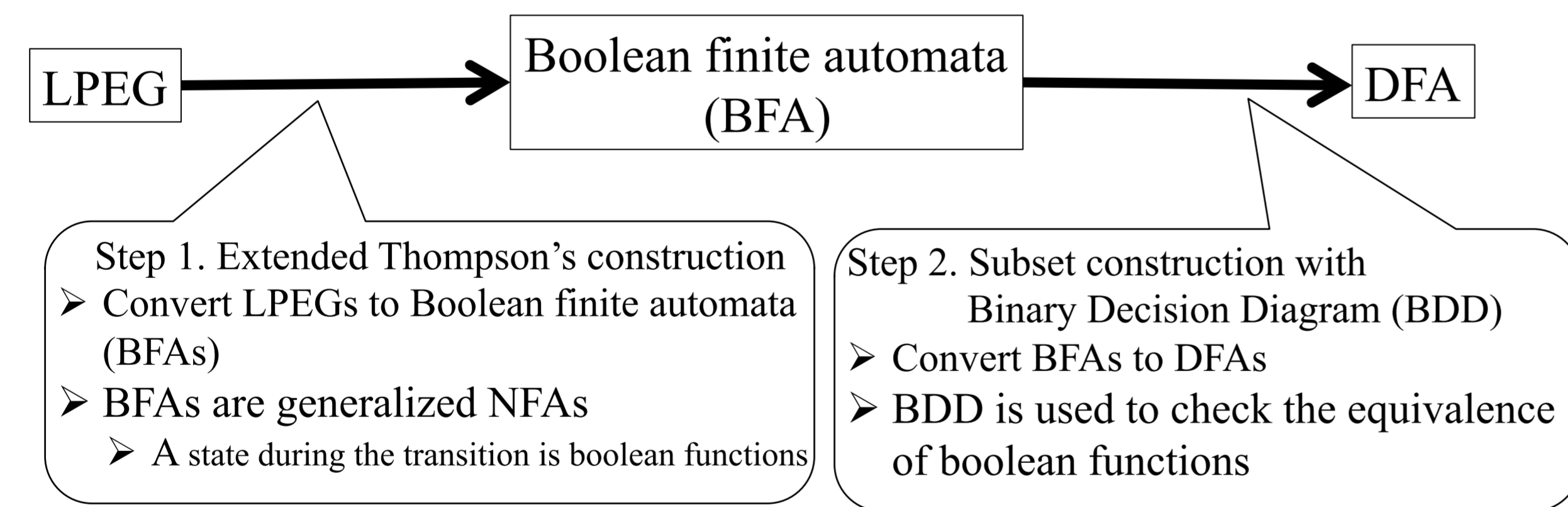
- Example 2 (not LPEG)

$A \leftarrow a A b / c A^* / d$  **bad**

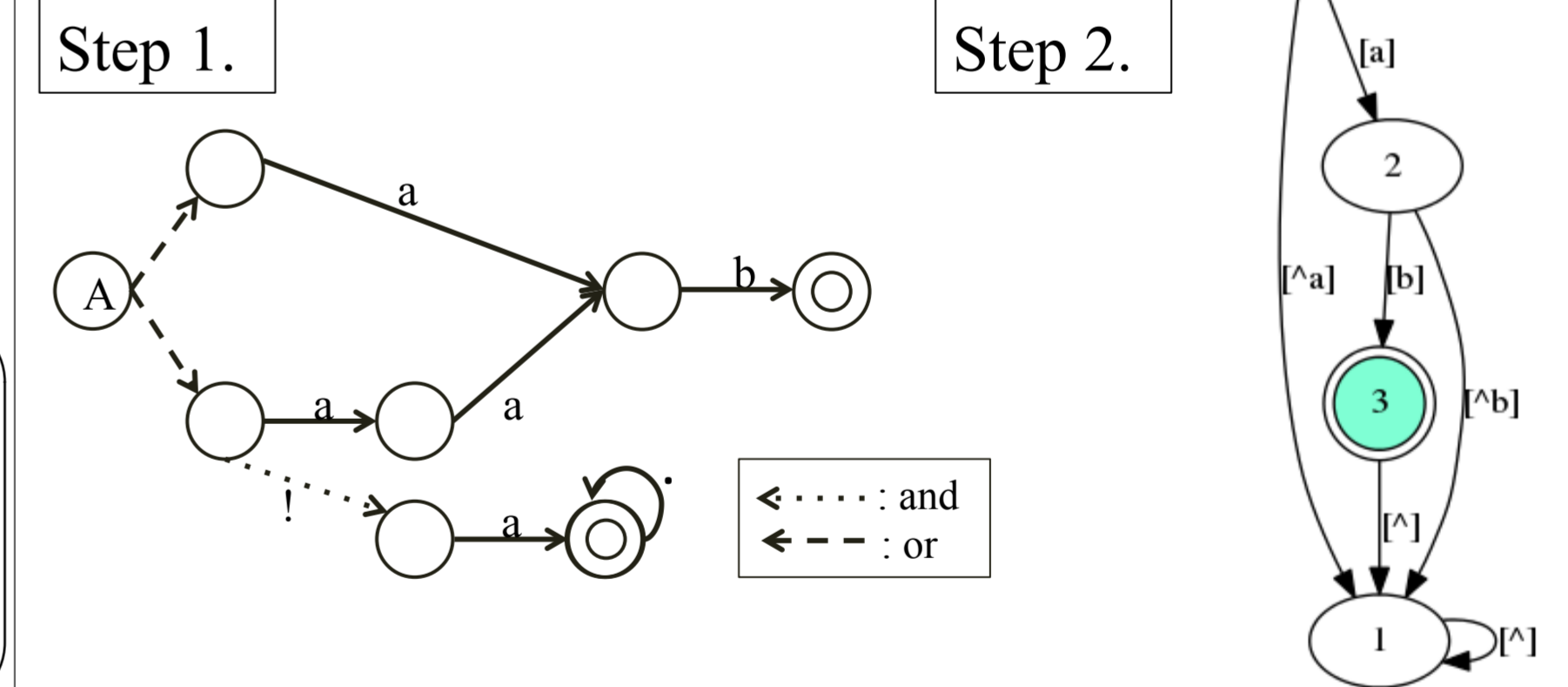
$e ::=$	$p$	
	$p A$	
	$p e$	
	$e/e$	
	$\&e e$	
	$!e e$	
$p ::=$	$\epsilon$	empty character
	$a$	any character
	$.$	any character
	$p p$	sequence
	$p/p$	prioritized choice
	$p^*$	zero or more repetition
	$\&p$	and-predicate
	$!p$	not-predicate

Figure 1. Syntax of a linear parsing expression

## Conversion from LPEGs to DFAs :



Example.  $A \leftarrow (a / aa) b$



## Performance : Applying the conversion to the parsers

We performed experiments to confirm the speed-up of PEG-based parsers by the partial conversion. The conversion is applied for each nonterminal iff the nonterminal is aLPEG. Figure 1 shows XML grammar in PEG.

Rules written in red letters are applied to the partial conversion.

### Experiment Method

We measured runtimes ten times in a row and calculated the averages of the runtimes other than the maximum runtime and the minimum runtime.

Table 2 and Table 3 show the result of the runtimes and the number of backtrackings, respectively.

Grammar	Input lines	Normal	DFA
XML	184,966	403ms	109ms
JSON	8,518	336ms	130ms

Grammar	Normal	DFA
XML	16,589,635	1,054,513
JSON	5,360,918	3,066,286

```
File <- PROLOG? DTD? Xml
Chunk <- Xml
Expr <- Xml
PROLOG <- '<?xml' (!'?'>')* '?'>' S*
DTD <- '<!>' (!'?'>')* '?'>' S*
Xml <- '<' Name S* Attribute* ('/>' / '>' S*
      (Content / COMMENT)* '</>' NAME '>' S*
Name <- NAME
NAME <- [A-Z_a-z:] ('-' / [0-9:A-Z_a-z])
Attribute <- Name S* '=' S* String S*
String <- '"' (!'"')* '"'
Content <- Xml / CDATAsec / Text
CDATAsec <- '<![CDATA[' CDATA '>]' S*
CDATA <- (!'>')* '<![CDATA[' CDATA '>]' CDATA
COMMENT <- '<!--' (!'-->')* '-->' S*
Text <- (!'<'>')*
S <- [ \t\r\n]
```

Figure 1. XML grammar in PEG