

# A century of parentheses languages

*with some amazing returns*

*ICTCS 2010, Camerino*

Stefano Crespi Reghizzi

Politecnico di Milano

# Parentheses in mathematical notation

Parentheses have appeared in algebraic writing in the XV-XVI century. Erasmus of Rotterdam calls them lunulae. Earlier and until the XVIII century, overline vinculum had been used for grouping literals into a term

$$\overline{aa + bb}^m \text{ instead of } (aa + bb)^m$$

Abstracting from the contents of parenthesized expressions, Walter von Dyck's (1856-1934) name has been given to the formal language every computer science student knows.

# Dyck's language

The alphabet has just two letters: the open/close parentheses  $(, )$  or *begin, end*, etc.

The language  $Cl(\varepsilon)$  is the equivalence class of all strings such that repeated deletions of well-parenthesized digram  $()$  reduce the string to the empty one  $\varepsilon$ .

$$(()())() \Rightarrow (())\emptyset() = (())() \Rightarrow (\emptyset)() \Rightarrow \emptyset() \Rightarrow \emptyset = \varepsilon$$

The equivalence class such that, after all deletions, the string, say,  $)$  is obtained is another formal language:  $Cl( )$  instead of  $Cl(\varepsilon)$ .

# Properties of Dyck's languages

Modest generalization: several matching pairs in the alphabet:

$$(, ), [, ], \{, \}, \dots$$

Obvious revision of cancellation rule.

Concatenating two or more times two such strings produces a Dyck string.

Reversing a Dyck string produces a Dyck language over the alphabet

$$([ ]) \xrightarrow{\text{reversal}} \overbrace{)}^{\text{open}} ] [ \overbrace{(}^{\text{close}}$$

Substituting a Dyck phrase for a character, say (, changes the equivalence class from  $Cl(\varepsilon)$  to  $Cl( ) )$  i.e., one closed paren in excess.

## Dyck's languages, grammars, push-downs and ...

Noam Chomsky's Context-Free grammar [1956] (or Bar-Hillel's Categorical g.) generates the Dyck language:

$S \rightarrow S S$       a phrase is the concatenation of 2 phrases

$S \rightarrow (S)$       a phrase is a parenthesized phrase

$S \rightarrow \varepsilon$       a phrase is the empty string

Word membership/parsing problem: given a string, is it a Dyck string? Deterministic push-down (LIFO) machine:

- Push on reading (
- Pop on reading ) and recognize if empty

Time complexity is linear (real-time).

# and ... queues

`\begin{parenthesis}`

equipped with a FIFO memory, a queue (or Post) machine recognizes the Anti-Dyck language [Vauquelin, Franchi-Zannettacci 1979], where “no parentheses match”.  
Cancellation rule:

$$\begin{array}{c} \text{contains no closing parens} \\ \underbrace{\left( \left( \left[ \right] \right) \right)}_{\underline{\quad}} \Rightarrow \left( \left[ \right] \right) \\ \underline{\quad} \end{array}$$
$$\left( \left[ \right] \right) \Rightarrow \left[ \right] \Rightarrow \varepsilon$$

Such languages are generated by breadth-first context-free grammars [Allewi, Cherubini, CR 1988].

`\end{parenthesis}`

# Parentheses unwelcome!

When parentheses proliferate they are hard to read. The number of parentheses can be reduced assigning precedences to operators:

$5 \times 3 + 8 \div 3 \times 9 + 7$  instead of  $(5 \times 3) + ((8 \div (3 \times 9)) + 7)$

$\times$  (and  $\div$ ) takes precedence over  $+$

$+$  yields precedence to  $\div$  (and  $\times$ )

$+$  yields to  $+$  (association from right to left)

$\times \succ +$

$+$   $\prec \div$

$+$   $\prec +$

Some people hate parentheses: Jan Lukasiewicz [1924] would write (without vincula) in reverse polish notation:

$\overline{5\ 3} \times \overline{\overline{8\ 3\ 9} \times} \div 7 + +$

## Floyd [1963]: operator precedence grammars

Idea: between all terminal characters there is a precedence relation:

- yields  $\langle$ , takes  $\rangle$
- equal-in-precedence,  $\doteq$ , between opening-closing pairs.

Compilers have extensively used Floyd grammars until the invention of deterministic methods (LL(k) Lewis et al. 1966, and LR(k) Knuth 1966);

Still popular for fast parsing of expressions.

Precedence relations are easily computed by grammar inspection.

# Parsing with precedence 1

Example: arithmetic expression with plus and times and with parens.

*Grammar* :  $E \rightarrow E + T \mid T, \quad T \rightarrow T \times F \mid F, \quad F \rightarrow (E) \mid \epsilon$

*Precedence matrix:*

	$a$	$+$	$\times$	$($	$)$
$a$		$\succ$	$\succ$		$\succ$
$+$	$\prec$	$\prec$	$\prec$	$\prec$	$\succ$
$\times$	$\prec$	$\prec$	$\succ$	$\prec$	$\succ$
$($	$\prec$	$\prec$	$\prec$	$\prec$	$\equiv$
$)$		$\succ$	$\succ$		$\succ$

## Parsing with precedence 2

Easy: syntax subtrees are delimited by  $\langle \dots \rangle$ , and may include subtrees separated by  $\dot{=}$ . No need to perform reductions from left to right.

$$\begin{array}{l}
 \vdash a \times ( a + a ) \vdash \\
 \vdash \langle a \rangle \times \langle ( \langle a \rangle + \langle a \rangle ) \rangle \vdash \\
 \vdash \langle a \rangle \times \langle ( \langle E + T \rangle ) \rangle \vdash \\
 \\
 \vdash \langle a \rangle \times \langle ( \dot{=} E ) \rangle \vdash \\
 \vdash \langle a \rangle \times F \rangle \vdash \\
 \dots \dots
 \end{array}$$

Clearly,  $\langle \dots \rangle$  act as parentheses.

I'll return to Floyd after various parenthesis models.

# Similarities with Regular languages REG

Few Properties of REG	preserved by	CF	Deterministic CF (DCF)
all Boolean		UNION	COMPL
Concatenation		YES	NO
Kleene Star		YES	NO
unique min. det. machine / grammar		NO	NO

Several scientists looked for “better” subfamilies of DCF:

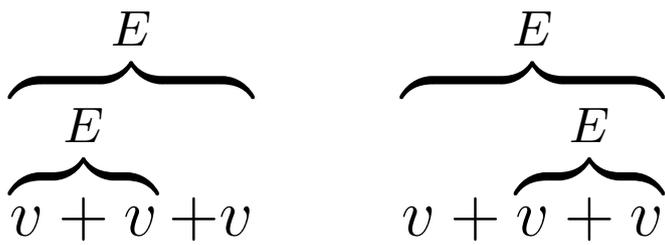
- Parenthesis grammars [McNaughton 1967, Knuth 1967], Tree automata [Thatcher 1967]
- Balanced Grammars [Berstel & Boasson 2002]
- Visibly Push-Down Automata VPD [Alur & Madhusudan 2004]

# Parentheses grammars and tree languages

Parentheses induce well-nested structures on strings.  
CF grammar rules are parenthesized [McNaughton]:

$$E \rightarrow E + E \mid v \quad \text{becomes} \quad E \rightarrow (E + E) \mid (v)$$

The ambiguous phrase  $v + v + v$  corresponds to different paren phrases



The diagram shows two parse trees for the expression  $v + v + v$ . In the first tree, the root  $E$  has two children:  $E$  and  $v$ . The left  $E$  child has two children:  $v$  and  $v$ . In the second tree, the root  $E$  has two children:  $v$  and  $E$ . The right  $E$  child has two children:  $v$  and  $v$ .

$$((v + v) + v) \quad (v + (v + v))$$

*Parentheses Languages* PL are DCF and math. similar to REG.

Very similar to *tree languages* [Thatcher 1967]

## REG-like properties of parentheses grammars

- Uniqueness of minimal grammar (in backwards-deterministic form).
- Grammars having the same set of rule patterns (*stencils*) define a Boolean algebra of languages.
- Non-Counting (aperiodic) REG languages [McNaughton, Papert, Schutzenberger] have counterparts within parentheses languages [CR, Guida, Mandrioli 1978]

Similar definitions and properties have been stated [Thomas] in the framework of tree automata.

## Further on: balanced strings and grammars 1

Surprisingly Dyck is not a parenthesis language:

$(())$   
missing external parens

$((()())\dots())$   
unbounded

PL are not closed under concatenation and Kleene star.  
[Knuth 1967] asked: is a given CF language a parens language? The answer involves an equivalent definition of well-parenthesizing for an alphabet including parens and possibly other “internal” letters.

A string is balanced if

- $\#$  open parens =  $\#$  of closed parens
- in every prefix,  $\#$  open parens  $\geq$   $\#$  closed parens

Dyck phrases are exactly the balanced strings.

# Letters associated to open paren

Every letter in a string is an associate of an open paren:

$$({}_0 c_1 ({}_2 c_3 )_4 ({}_5 )_6 )_7$$

letter 1 is associate of  $({}_0$ , letter 3 of  $({}_2$ , letter 4 of  $({}_0$   
[Knuth] A CF language is a parens language iff

- every phrase is balanced and
- every open parens has bounded number of associates.

$$S \rightarrow XY \quad X \rightarrow ( )c)X) \mid (d \quad Y \rightarrow (Y(c)) \mid e)$$

is a parens language, though grammar is not parenthesized.

## Paren nesting in human and artificial languages

- Natural languages rarely exhibit deeply nested structures
- although in principle they are possible  
*der Mann der die Frau die das Kind das die Katze füttert sieht liebt schläft*
- inner clauses are rarely marked by parens or by words acting as opening / closing tags
- good writers moderately use *parentheticals*, because they depart from the main subject

# Parens in computer languages

- All programming languages have parenthetical constructs, perhaps exaggeratedly so in Algol 68

*begin ... end, do ... od, if ... fi, case ... esac*

- Mark-up or semi-structured (web) documents (e.g. XML) are deeply and widely nested; visible open/close tags delimit structures:

```
<div id="accessorapido">
  <ul>
    <li><a href="#barranavigazione">Navigaz
    <li><a href="#avvisi">Avvisi</a></li>
    <li><a href="#contenutoprincipale">cont
    <li><a href="#barrainformazioniaggiunti
  </ul>
</div>
```

# Balanced grammars [Berstel & Boasson 2000]

Allow RegExpr in righthand sides of rules:

$$S \rightarrow (Y^*) \quad Y \rightarrow []$$

phrases:  $()$ ,  $([])$ ,  $([] [])$ ,  $\dots$

Several properties of regular languages hold for balanced languages:

Boolean closures, uniqueness of minimal grammar.

An elegant formalism for XML-like languages.

# Visibly PushDown languages [Alur&Madhusudan 2004]

Restricted type of deterministic pushdown machine.

Motivated by:

- model-checking of programs (i.e.  $\infty$  state systems)
- XML

VPD Push-Down Machine:

- pushes a call (=open) letter, changing state
- pops on a return (=close) letter, if a matching call letter is on top, changing state
- changes state on a return letter, if stack empty
- on an internal letter, changes state without using stack
- accepts by final state.

# VPD versus balanced lang. and REG

- 3-partite alphabet:  $\overbrace{\Sigma_{call}}^{\text{open parens}} \cup \overbrace{\Sigma_{return}}^{\text{closed parens}} \cup \overbrace{\Sigma_{internal}}^{\text{others}}$
- no bijection open-close: e.g.,  $( ]$  and  $( )$
- unbalanced returns may occur as prefix of a word, and unbalanced calls as suffix

$\overbrace{r_1 r_2}^{\text{unbal. returns}}$ 
 $\overbrace{c_1 c_1 \dots r_1 r_1}^{\text{balanced}}$ 
 $\underbrace{i_1}_{\text{internal}}$ 
 $\overbrace{c_1 c_1 \dots \underbrace{i_2}_{\text{internal}} r_2 \dots \underbrace{i_2}_{\text{internal}} r_2}^{\text{balanced}}$ 
 $\overbrace{c_2 \dots c_2}^{\text{unbal. calls}}$

- $\text{REGULAR} \subset \text{BALANCED} \subset \text{VPD} \subset \text{Deterministic CF}$
- determinization, minimization, uniqueness
- closure and decidability properties  $\approx$  to REGULAR
- real-time deterministic parsing

# Examples

- Dyck equivalence class ‘)’:

$(( ( ( ) ) ) )$  parsed as  $\overbrace{(( ( ( ) ) ) )}$  not as  $( \overbrace{(( ( ) ) )} )$

- Program execution modelled by VPD, e.g., abnormal termination of procedure call:

$C_{main} C_A \overbrace{C_B \dot{i} instruction \dots r_B C_B \dot{i} instr. \dot{i} except.}$

# Plenty of research on VPD

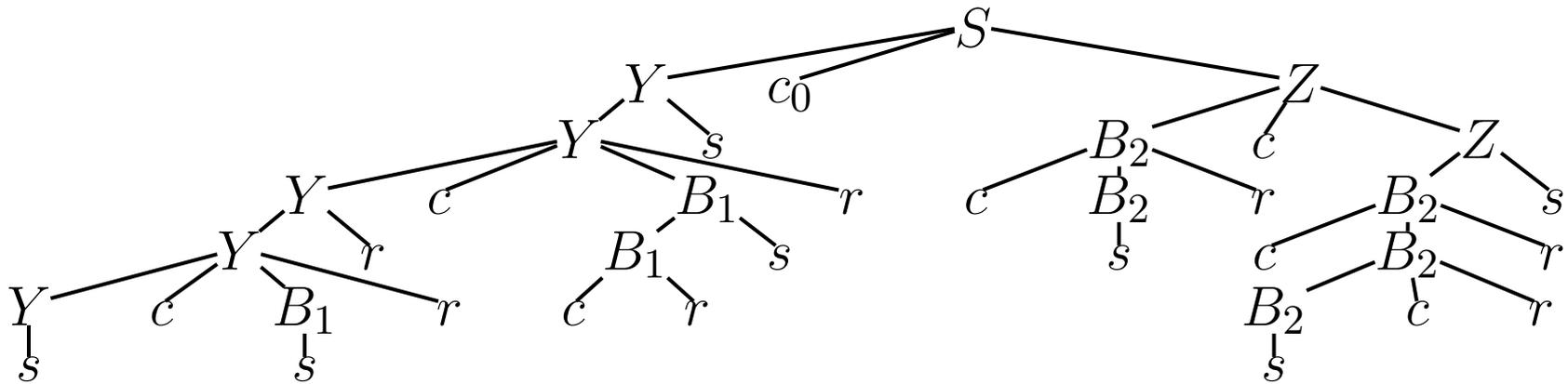
- Monadic Second Order Logic of VPD languages
- Grammatical formulations of VPD have been provided
- $\omega$  (infinitary) languages
- Bisimulation equivalence
- Decidable problems
- VPD games
- Complexity of membership problem w.r.t. grammar size [La Torre, Napoli, Parente 2006]
- Comparison with synchronized [Caucal] and height-deterministic [Nowotka & Srba] languages

# But are VPD languages that new?

- Floyd [1963] Operator Precedence [1963] Languages  
FL strictly include VPD [CR & Mandrioli 2009]
- The precedence relations of VPD languages have a particular form
- A FL grammar with such form of precedences generates a VPD lang.
- The same formal properties hold for FL and VPD
  - Boolean closure [CR, Mandrioli, Martin 1978]
  - Reversal
  - Closures under concatenation, star, suffix, prefix [CR & Mandrioli 2010]

# Precedence relations of VPD

Structure of FL grammar of a VPD lang. :



Precedence relations:

	$c \in \Sigma_{call}$	$r \in \Sigma_{return}$	$s \in \Sigma_{internal}$
$\Sigma_{call}$	$\triangleleft$	$\dot{=}$	$\triangleleft$
$\Sigma_{return}$	$\triangleright$	$\triangleright$	$\triangleright$
$\Sigma_{internal}$	$\triangleright$	$\triangleright$	$\triangleright$

A language is VPD iff it is generated by a FL grammar with such precedences

# Limitations of VPD

- Open / close tags must be letter-disjoint

$a \overbrace{a \dots b} b$        $c \overbrace{c \dots d} d$

contradict

$e \overbrace{e \dots ac} ac$

- fixed syntax structure in many cases not structurally adequate

bad  $3 \overbrace{+ 5} \times 7$       good  $3 \overbrace{+ 5} \times 7$

Both cases are correctly handled by Floyd grammars.

# Floyd lang. as generalized parens lang.

- Functional notation uses nested parens:  
 $ADD(a, MULT(c, ADD(d, e)))$
- more readable in *infix notation* with precedences:  
a ADD b MULT ( c ADD d)
- For ternary operators  
 $IF\_THEN\_ELSE(c_1, s_2, s_3)$   
*mixfix notation* imitates natural language

$$\underline{IF} \ c_1 \ \underline{THEN} \ s_2 \ \underline{ELSE} \ s_3$$

with precedences:

- $\underline{IF} \doteq \underline{THEN} \doteq \underline{ELSE}$
- $\underline{IF} \triangleleft$  1st symbol of  $c_1$
- last symbol of  $c_1 \triangleright \underline{THEN}$

# Conclusione semiseria

[dal *Breve glossario di retorica e metrica*]  
Parentesi o frase incidentale è l'aggiunta di elementi non necessari o di precisazioni all'interno di una frase. È segnalata dalle parentesi o dalle virgole.

Sono contrito di avervi intrattenuto per un'ora parlandovi di elementi non necessari!  
O forse l'informatica teorica è la scienza del non-necessario  
... ma illuminante?