

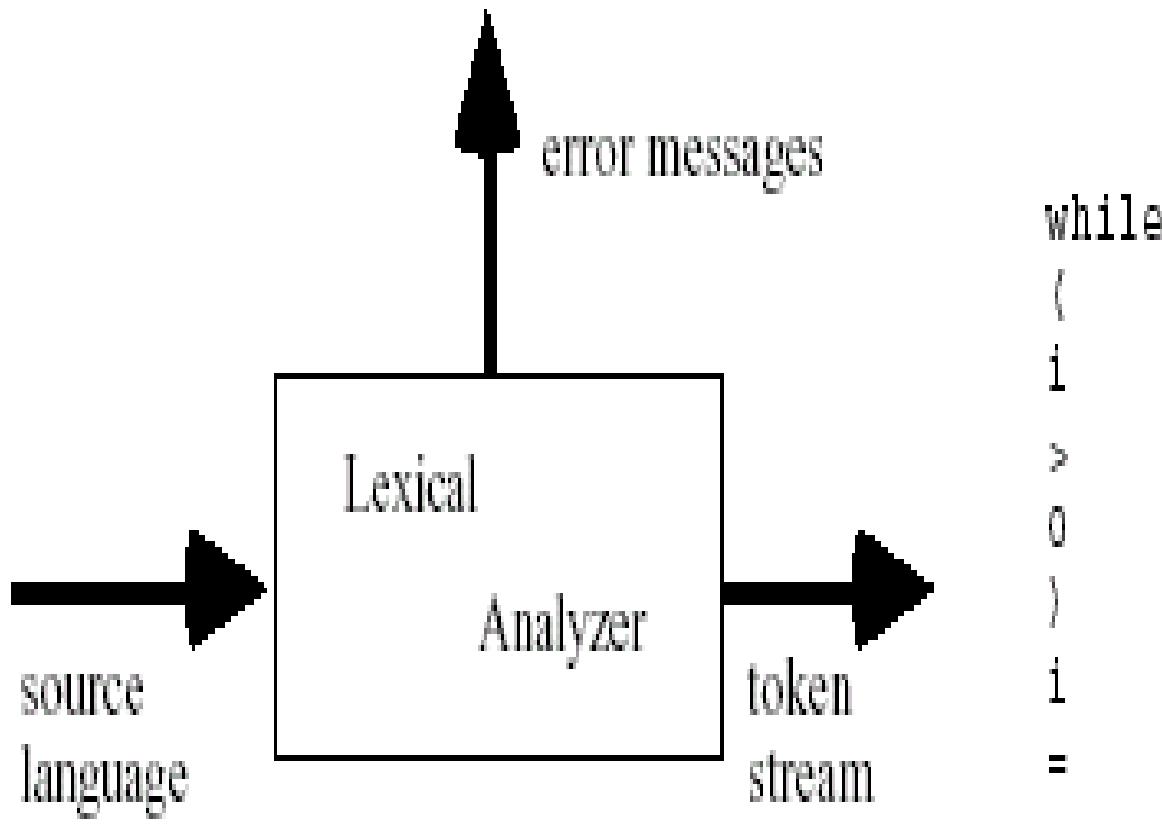
# **The Scanner (Lexical Analyzer)**

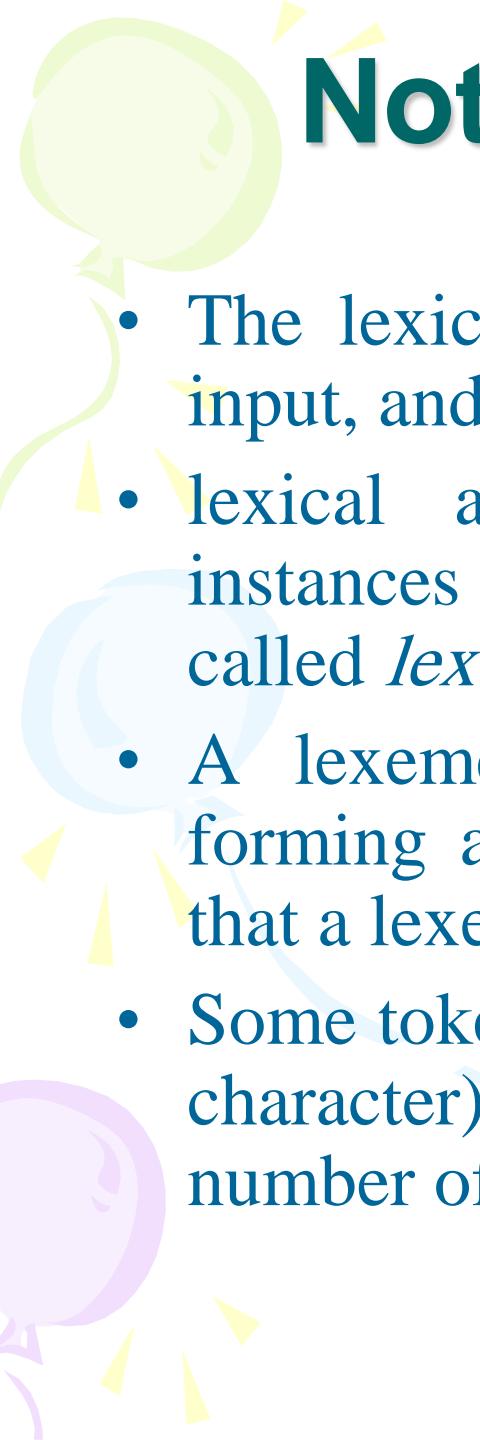
# The Scanner (Lexical Analyzer)

- **The basics**

- *Lexical analysis* or *scanning* is the process where the stream of characters making up the source program is read from left-to-right and grouped into tokens.
- *Tokens* are sequences of characters with a collective meaning.
- There are usually only a small number of tokens for a programming language: constants (integer, double, char, string, etc.), operators (arithmetic, relational, logical), punctuation, and reserved words.

# lexical analysis





# Notes on lexical analysis

- The lexical analyzer takes a source program as input, and produces a stream of tokens as output.
- Lexical analyzer might recognize particular instances of tokens, such specific instances are called *lexemes*.
- A lexeme is the actual character sequence forming a token, the token is the general class that a lexeme belongs to.
- Some tokens have exactly one lexeme (e.g. the > character); for others, there are an infinite number of lexemes (e.g. integer constants).

# Notes on lexical analysis

- There are two primary methods for implementing a scanner.
- The first is a program that is hard coded to perform the scanning tasks. It is .. **Scanner implementation #1: loop & switch**
- The second is the use of regular expressions and finite automata to model the scanning process. It is..**Scanner implementation #2: regular exp-ressions & finite automata.**

# lexical analysis

## • التعرف على الـ *tokens*

خلال ذلك سيتم الاعتماد على هذا النحو كمثال

## ***tokens***

للشرح .

## • مثال: افترض هذا الجزء من النحو :

• Stmt  $\rightarrow$  if expr **then** stmt

| if expr **then** stmt **else** stmt

|  $\epsilon$

• expr  $\rightarrow$  term **relop** term

| term

• term  $\rightarrow$  id | num

# lexical analysis

- حيث النهايات (terminals)  
**if , then , else , relop , id , num . . .**

- مكونة من فئات من التسلسلات معطاه بواسطة التعريفات المنتظمة الآتية:

• **if** → if

• **then** → then

• **else** → else

• **relop** → < |<= | > | >= | <> | =

• **id** → letter ( letter | digit )\*

• **num** → digit+ (.Digit+ )? (E (+|-)? digit+ )?

• **letter** → A|B....|Z|a|b|....|z |A – za – z|

• **digit** → 0|1|.....|9 → |0 -9 |

# lexical analysis

- وفي هذا الجزء من اللغة فإن المحلل الإملائي Lexical analyzer يقوم بالتعرف على الكلمات الحاكمة if , then , else relop , id ،-- lexeme بالإضافة لل keywords num
- ولاختصار الحالة سنفترض:
  - \* أن الكلمات الحاكمة keywords كلمات ممحوزة أي لا يمكن استخدامها ك identifier
  - . real num \*\* سيمثل عدد بدون اشارة أو عدد حقيقي
  - هذا بالإضافة إلى إننا سنفترض lexemes سيتم فصلها بواسطة White Space
  - nonnull sequence of newline , tabs (ws) blanks
  - وسيقوم المحلل الإملائي بفصل white space وسيقوم بذلك بمقارنة التتابع بالتعريف المنتظم

# Lexical analysis review

## ملخص للمحلل اللغوي

- Lexical analysis: the process by which a sequence of characters (input stream) is broken up into its elementary words (tokens)
  - أنه العملية التي بها يتم تقطيع التسلسلات الدالة إلى الكلمات الأولية
- Our goal: come up with an algorithm to scan an input stream and identify the tokens
  - إن الهدف هو عمل خوارزم يقوم بمسح التسلسلات الدالة و التعرف للمعرفات tokens
- What we know: a general description of what we know about tokens
  - ما نعرفه هو توصيف عام للعلامات the tokens

# Lexical analysis review

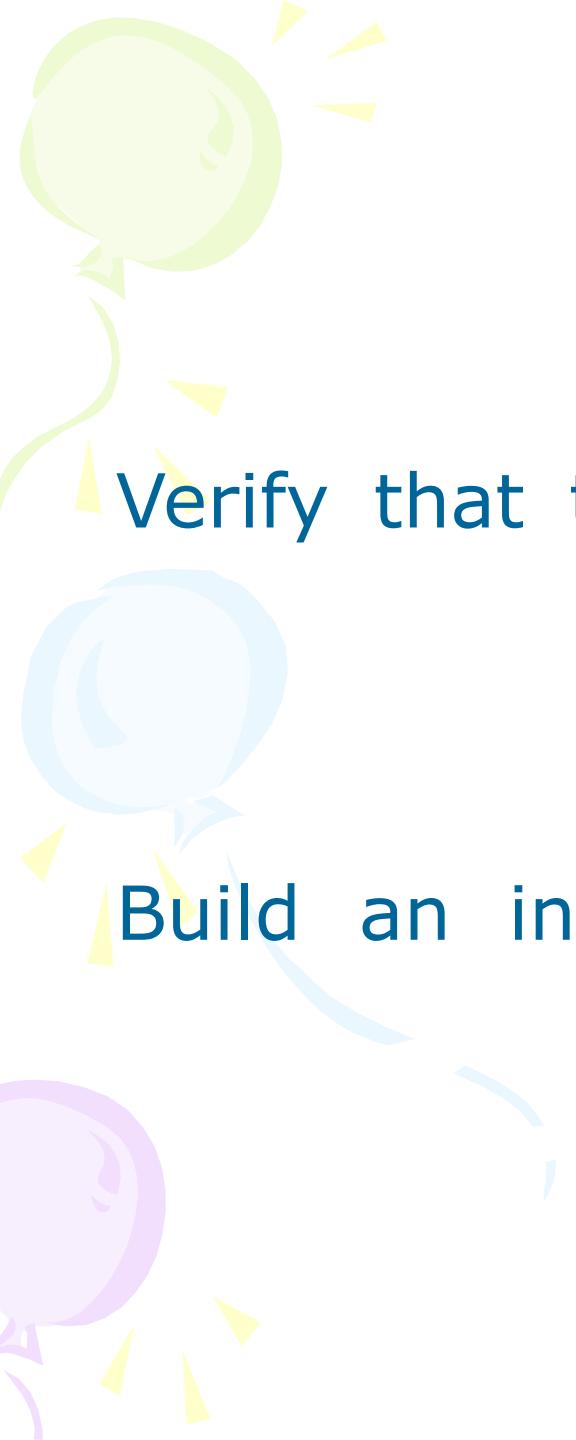
## ملخص للمحال اللغوي

- Step 1: Use regular expressions to formally describe the tokens
  - استخدم التعبير المنتظم في الوصف المحدد للعلامة token
- Step 2: Build a finite automaton out of the regular expressions, i.e. a machine (algorithm) that can recognize words described by those regular expressions.
  - قم ببناء الآلة ذاتية الحركة من التعبير المنتظم أي خوارزم الآلة التي تعرف على الكلمات الموصفة بذلك التعبير المنتظم
- Step 3: Add some additional functionality: error handling, symbol table access.
  - أضف بعض الدوال الازمة لذلك الحل مثل كيفية التعامل مع الأخطاء على جدول الرموز

# Lexical analysis review

## ملخص للمحلل اللغوي

- Result: a program (machine, algorithm) that can read the input stream and return tokens.
- النتيجة: المحلل الإملائي هو عبارة عن برنامج <آلـةـ خوارزم> والذي يستطيع قراءة التتابع الداخـل حيث يقوم بإعطاء أو أخـرـاج علامـات Tokens
- All this can be done automatically by a scanner generator (such as Lex) as long as we provide the regular definitions and specify the “additional” actions.
- كل هذا يمكن عمله آلـياـ باستخدام ”المولد الماسح“ مثل (Lex) طالما تم عمل التعريفات المنتظمة وأيضا إضافة بعض العمليـات الإضافـية



# Syntax analysis

## المحل النحوي

During this phase:

Verify that the program is syntactically correct

• خلال تلك المرحلة يتم الآتي :-

• التحقق من أن البرنامج صحيح اعرابيا

Build an internal representation of the source code

• \*بناء تمثيل داخلي لبرنامج (كود) المنهج

Handle errors

• \*إعطاء وطباعة الأخطاء

# Context-Free Grammars (CFGs)

- A scanner recognizes words of the language
- إن الماسح (المحلل الإملائي) يتعرف على كلمات اللغة

- A parser recognizes sequences (strings) of tokens.

- بينما المعرف يتعرف على تتابعات العلامات Tokens مثل :-
- e.g. **(id+num)\*id**

- is a sequence of tokens that "builds" an expression.

- -هذا التعبير عبارة عن تسلسلات من العلامات والتي من خلالها يتم بناء تعبير

- We need a way to formally describe such patterns of tokens.

- لذا يلزمنا لعمل المحلل الإعرابي طريقة معروفة لتوصيف تلك النماذج للTokens

# Context-Free Grammars (CFGs)

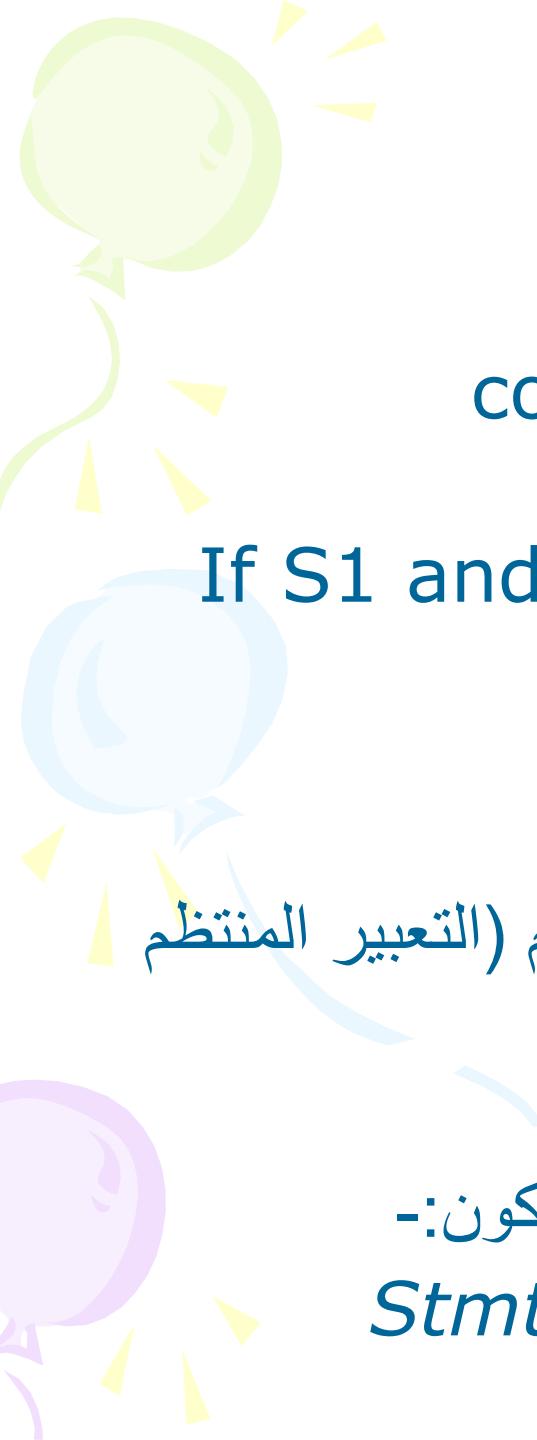
Regular expressions are not powerful enough (for example, they cannot be used to describe balanced parentheses or block structure)

إن التعبيرات المنتظمة ليست بالقوة الكافية لعمل المعرف -فمثلا لا تستطيع عمل الأقواس وإتزانها- أي الفتح والقفل للأقواس

Context-free grammars, however, are!

بينما تستطيع تقنية CFG عمل ذلك

كذلك معظم لغات البرمجة بناؤها يحتوي على تكرير للبناء مما يمكن تعريفها بواسطة CFG

- 
- مثال إذا كان عندنا تعبير مشروط conditional statement معرفة بالقاعدة الآتية :

If S1 and S2 are statements and E is an expression , then :

- “ If E then S1 else S2” is a statement. (1)

• وهذا التعبير لا يمكن وضعه في صورة تعبير منتظم (التعبير المنتظم يتم استخدامه في البناء الإملائي كما درسنا)

• لذا فإنه يمكن استخدام نحو المنتجات

• في (1) ليكون:- grammar production

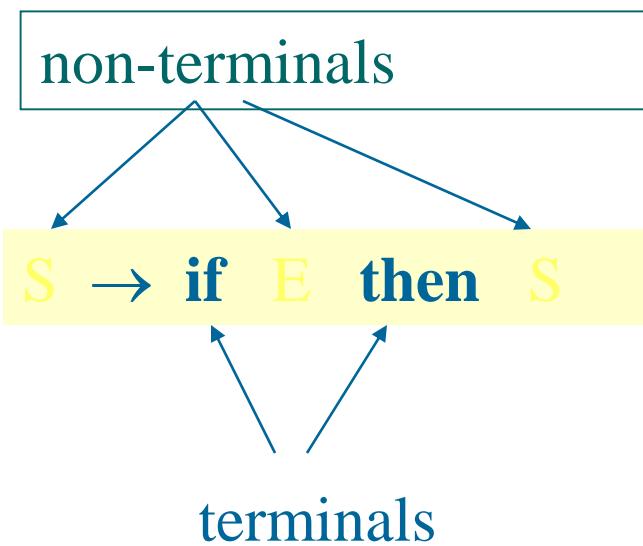
•  $\text{Stmt} \rightarrow \text{if expr } \underline{\text{then}} \text{ stmt } \underline{\text{else}} \text{ stmt}$

# Context-Free Grammars (CFGs)

- A CFG is a language generator that uses "rules" to generate valid strings of a language.
- إن تقنية CFG عبارة عن مولد للغة والذي يستخدم قواعد لتوليد تتابعات محققة من اللغة
- Intermediate symbols are called **non-terminals** الرموز المتوسطة تسمى نهايات غير طرفية
- Atomic symbols are called **terminals** الرموز الأساسية تسمى نهايات طرفية

# Context-Free Grammars (CFGs)

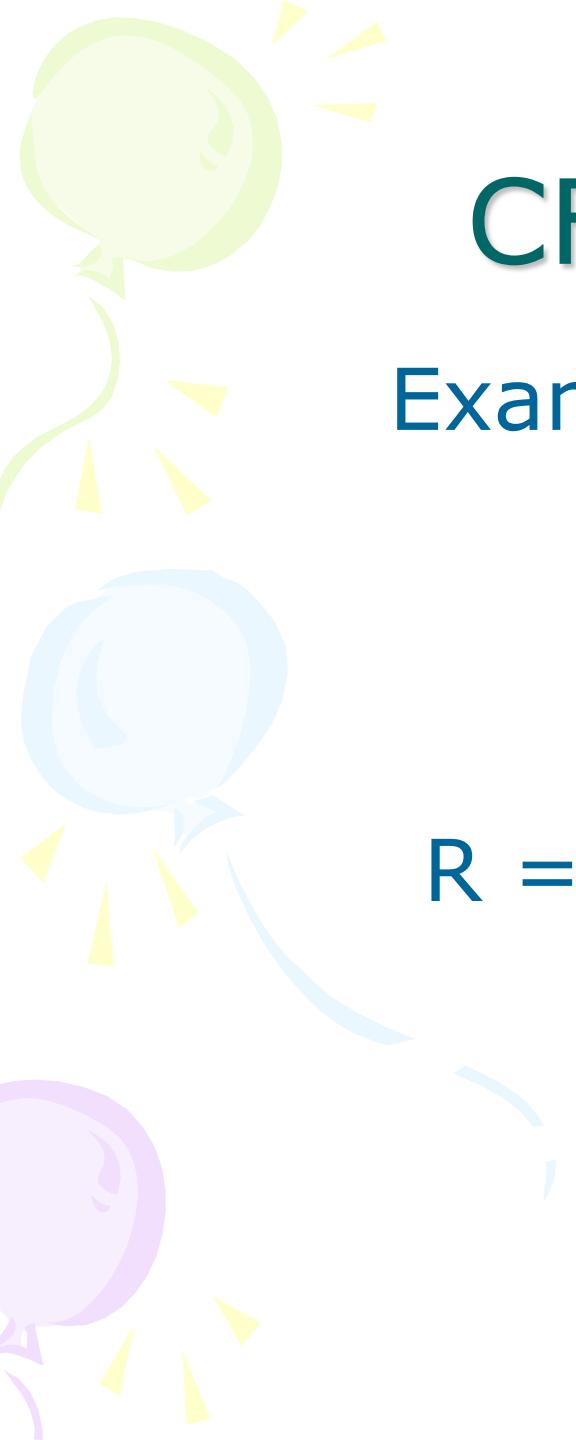
Example: •





# CFGs : definition

- A CFG is a quadruple  $(V, \Sigma, R, S)$  • where
  - $V$  is an alphabet consisting of – terminals and nonterminals
  - $\Sigma \subseteq V$  is the set of terminals –
  - $S \in V - \Sigma$  is the starting symbol –
  - $R \subseteq (V - \Sigma) \times V^*$  is a set of rules –



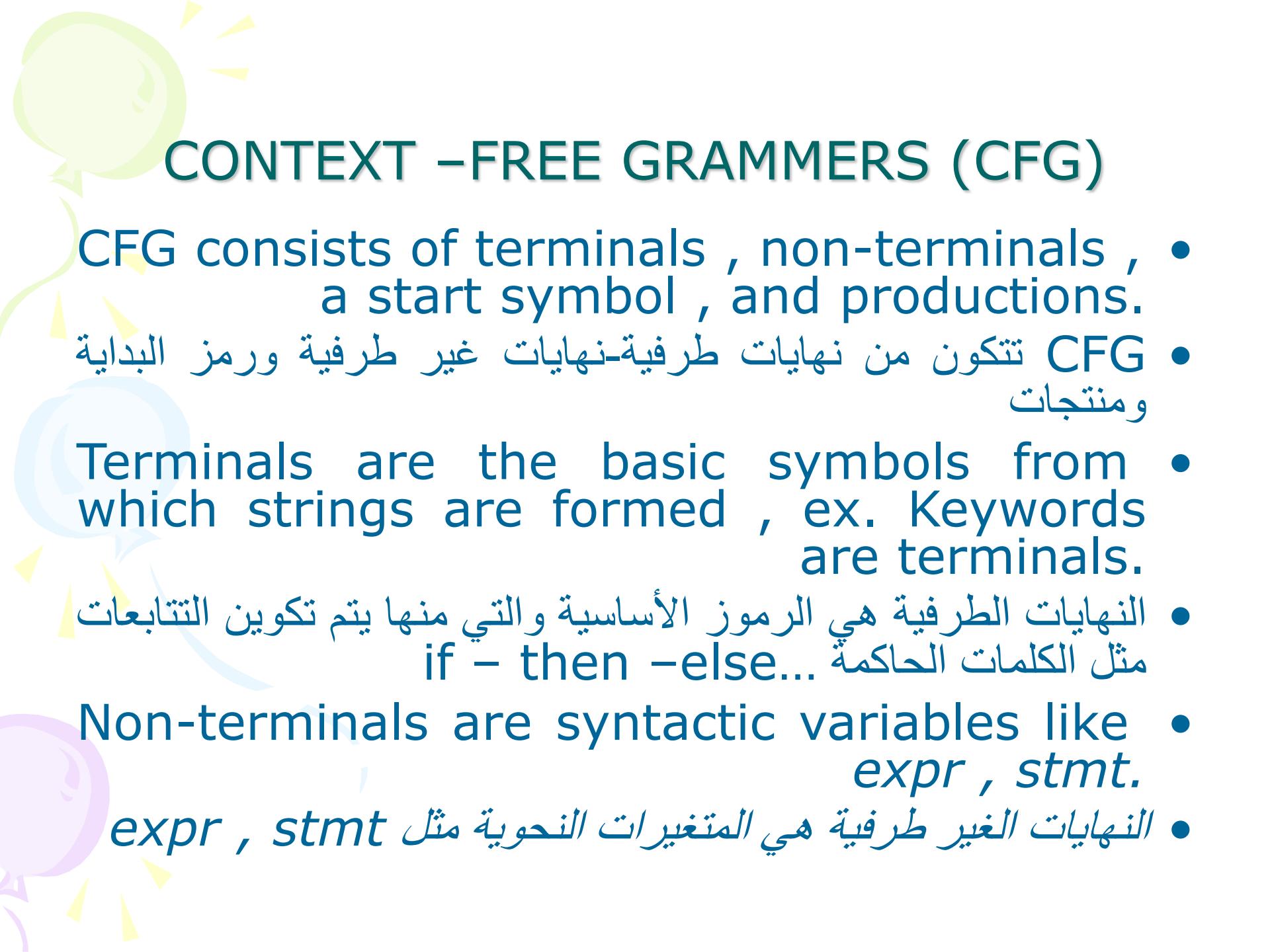
# CFGs : definition

Example: Consider  $(V, \Sigma, R, S)$  •  
where

$$V = \{S, (, )\} -$$

$$\Sigma = \{ (, ) \} -$$

$$R = \{ S \rightarrow (S), S \rightarrow SS, S \rightarrow \varepsilon \} -$$



# CONTEXT -FREE GRAMMERS (CFG)

- CFG consists of terminals , non-terminals , a start symbol , and productions.
- CFG تكون من نهايات طرفية-نهايات غير طرفية ورمز البداية ومنتجات
- Terminals are the basic symbols from which strings are formed , ex. Keywords are terminals.
- النهايات الطرفية هي الرموز الأساسية والتي منها يتم تكوين التتابعات مثل الكلمات الحاكمة if - then -else...
- Non-terminals are syntactic variables like *expr* , *stmt*.
- النهايات الغير طرفية هي المتغيرات النحوية مثل *expr* , *stmt*

# CONTEXT -FREE GRAMMERS (CFG)

In a grammar , one non-terminal is distinguished as start symbol and the set of strings it denotes is the language defined by the grammar.

وفي النحو فإن أحد هذه النهايات الغير طرفية يتم اختياره كرمز البداية ثم التتابعات التي يشير إليها هي اللغة التي تخضع لائل القواعد

The production of a grammar specify the manner in which the terminals and non-terminals can be combined

# CONTEXT -FREE GRAMMERS (CFG)

- ومنتجات النحو عبارة عن الطريقة التي بها يتم ربط النهايات الطرفية وغير الطرفية لتكوين تتابعات النحو
- Each production consists of a non-terminal, followed by an arrow (sometimes the symbol ::= is used in place of the arrow) , followed by a string of non-terminals and terminals.
- حيث أن كل منتج عبارة عن نهاية غير طرفية متبوعة بسهم أو = أو ::= متبوعا بمتسلسلات من النهايات الطرفية وغير الطرفية

# CONTEXT -FREE GRAMMERS (CFG)

$Expr \rightarrow (expr)$

$Expr \rightarrow id$

$Op \rightarrow -$

$Op \rightarrow /$

In this grammar , the terminal symbols are :-

$id , + , - , * , ↑, ( , )$

The non-terminals are  $expr$  and  $op$

Ex. •

$expr \rightarrow expr op expr$  •

$Expr \rightarrow - expr$  •

$Op \rightarrow +$  •

$Op \rightarrow *$  •

$Op \rightarrow /$  •

# Notational Conventions

## مصطلاحات رمزية

- ولتفادي ذكر دائماً أنّ هذا نهاية طرفية terminal وهذا نهاية غير طرفية non-terminal. فأننا سنقوم بعمل اصطلاحات رمزية متعارف عليها وذلك خلال ما تبقى من دراسة هذه المادة

### 1-These symbols are terminals:

- هذه الرموز نهاية طرفية
- i) Lower -case letters early in the alphabet such as a, b, c... الحروف الأنجلizية الصغيرة والتي في البداية مثل علامات
- ii) Operator symbols such as +, -, \*, ..... العمليات الرياضية
- iii) Punctuation symbols such as parentheses (,), comma, ..etc الرموز التي تستخدم بين الكلمات والجمل مثل الأعداد
- iv) The digits 0 ,1 ,2 ,3 ,4 ,.....,9
- Boldface strings such as **id** or **if** (For us may be id or if)

# Notational Conventions

## مصطاحات رمزية

- **2-These symbols are non- terminals:**
- i)Upper - case letters early in the alphabet such as A , B, C, ...
- ii)The letter S , which, when it appears , is usually the start symbol.
- iii)Lower -case italic names such as *expr* or *stmt*
- 3-Uppercase letters late in the alphabet , such as X ,Y,Z represent grammar symbols , that is , either non-terminals or terminals

# Notational Conventions

## مصطلاحات رمزية

- 4-Lower - case letters late in the alphabet u ,v , ....z , represent strings of terminals
- 5- Lower - case Greek letters ,  $\alpha$  ,  $\beta$  ,  $\gamma$ , for example , represents strings of grammar symbols , i.e:
  - معناها  $A \rightarrow \alpha$  ,
- If  $A \rightarrow \alpha_1$  ,  $A \rightarrow \alpha_2$  ,  $A \rightarrow \alpha_3$  are all production with A on the left ( we call them A-productions) we may write
- $A \rightarrow \alpha_1 | \alpha_2 | \alpha_3.....| \alpha_k$ . We call  $\alpha_1, \alpha_2, \alpha_3, ..., \alpha_k$  the alternatives for A
- Unless otherwise stated , the left side of the first production is the start symbol

# Notational Conventions

## مصطلاحات رمزية

- Example :-

- وبإستخدام هذه المصطلحات سيتم التعبير عن النحو الآتي وكما هو موضح

$expr \rightarrow expr \ op \ expr$

$xpr \rightarrow - \ expr$

$Op \rightarrow \uparrow$

- $Expr \rightarrow (expr)$

$Op \rightarrow +$

$Op \rightarrow -$

- $Op \rightarrow /$

$E \rightarrow E \ A \ E \ | \ (E) \ | \ -E \ | \ id$

$A \rightarrow + \ | \ - \ | \ * \ | \ ^$

- وهكذا فإن المصطلحات الرمزية تخبرنا أن  $E, A$  terminals

- And  $E$  start symbol . The remaining

# Derivations

## إثباتات

- ولإثبات أن النحو يتم بواسطته تعريف اللغة
  - وكمثال على ذلك فإن هذا النحو والذي يتم استخدامه للتعبير عن التعبيرات الرياضية حيث  $E$  تمثل تعبير  $E \rightarrow E+E \mid E * E \mid (E) \mid -E \mid id$
  - $E \rightarrow -E$
  - معناها أن تعبير مسبوقا بعلامة  $-$  أيضاً تعبير ويمكن اختصاره كالتالي:-  
 $E \Rightarrow -E$  Which means " $E$  derives  $-E$ "  
أي ينتج عنها  
وهكذا فإن
  - $E \Rightarrow (E)$
  - Which means we could also replace one instance of an  $E$  in any string of grammar symbols by  $(E)$  as:-

# Derivations

## إثباتات

النحو  $E \rightarrow E+E \mid E*E \mid (E) \mid -E \mid id$

المستخدم

$E \Rightarrow E * E \Rightarrow E * (E)$

$E \Rightarrow E * E \Rightarrow (E) * E$

وهكذا فيمكن أخذ  $E$  واحدة ثم تكرارها وموضعين منتجات بأي ترتيب وذلك لعمل تتابع ما

- $(id)$

$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(id)$

$E \Rightarrow -(id)$  من

# Derivations

## إثباتات

النحو  $E \rightarrow E+E \mid E^*E \mid (E) \mid -E \mid id$

المستخدم

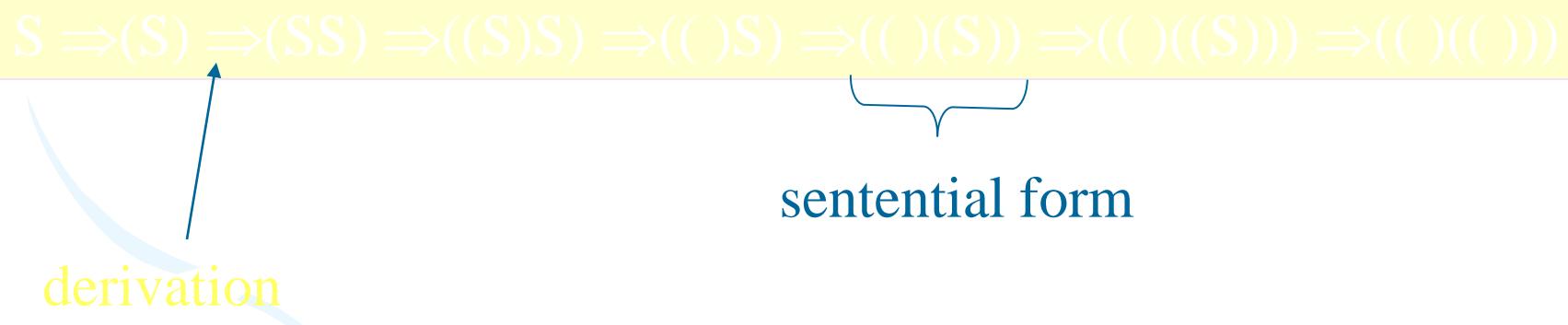
• مثال التسلسل (**id+id**) - عبارة عن جملة في النحو السابق

• هذا لأنه يمكن إيجاد إثبات كالتالي

$$\begin{aligned} E &\Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow \dots \Rightarrow -(id+E) \Rightarrow - \\ &\quad (id+id) \end{aligned}$$

# CFGs : derivations

Example: Consider  $(V, \Sigma, R, S)$  where •  
 $\Sigma = \{ (, ) \} -$        $V = \{S, (, )\} -$   
 $R = \{ S \rightarrow (S), S \rightarrow SS, S \rightarrow \varepsilon \} -$



This example demonstrates a **leftmost derivation** : one where we always expand the leftmost non-terminal in the sentential form.