

## ΚΕΦΑΛΑΙΟ 3 ΣΥΝΤΑΚΤΙΚΑ ΣΤΟΙΧΕΙΑ ΓΛΩΣΣΩΝ ΠΡΟΓΡΑΜΜΑΤΙΣΜΟΥ

### Στόχος

Στόχος του Κεφαλαίου αυτού είναι να περιγράψει με σύντομο τρόπο κάποια βασικά στοιχεία γλωσσών προγραμματισμού τύπου Context-free (Γραμματικές χωρίς συμφραζόμενα) μια και στη συνέχεια του βιβλίου θα αναφερόμαστε μόνο σε τέτοιες γλώσσες.

### Προσδοκώμενα αποτελέσματα

Όταν θα έχετε μελετήσει το κεφάλαιο αυτό θα μπορείτε να:

- εξηγήσετε γιατί χρησιμοποιούμε Context-free γραμματικές στην σύνταξη των γλωσσών προγραμματισμού
- αιτιολογήσετε γιατί οι Context-free γραμματικές είναι ισχυρότερες από τις Κανονικές Εκφράσεις ή τις Κανονικές γραμματικές
- δείξετε ότι μια γραμματική είναι διαφορούμενη
- αποδιοφοροποιήσετε μια διαφορούμενη γραμματική
- περιγράψετε πώς παράγονται τα δένδρα ανίχνευσης ενός Συντακτικού Αναλυτή
- δώσετε τους τυπικούς ορισμούς του Chomsky για γραμματικές τύπου 0, 1,2, και 3.

### Έννοιες κλειδιά

Context-free γραμματικές, γραμματικές Chomsky, διαφορούμενες γραμματικές, συντακτικές κατηγορίες.

### Εισαγωγικές παρατηρήσεις

Είδαμε στο προηγούμενο κεφάλαιο ότι η λεκτική δομή των tokens μπορεί να προσδιοριστεί με Κανονικές Εκφράσεις (ΚΕ) και ότι από τις ΚΕ μπορούμε αυτόματα να κατασκευάσουμε ένα Λεκτικό Αναλυτή που αναγνωρίζει τις λεκτικές δομές (tokens) μιας γλώσσας.

Με παρόμοιο τρόπο θα χειριστούμε στα επόμενα κεφάλαια την Συντακτική Ανάλυση (ΣΑ). Για τον συντακτικό καθορισμό μιας γλώσσας προγραμματισμού θα χρησιμοποιήσουμε ένα συμβολισμό που καλείται BNF (Backus-Naur-Form), και θα ασχοληθούμε με γραμματικές τύπου Context-Free ή αλλιώς γραμματικές Chomsky τύπου 2 (τις γραμματικές αυτές θα τις ορίσουμε στην ενότητα 3.5).

Θα περιγράψουμε πώς μια γραμματική ορίζει μια γλώσσα, και θα εξετάσουμε ποιά χαρακτηριστικά γλωσσών προγραμματισμού, μπορούν ή δεν μπορούν να οριστούν από γραμματικές τύπου Context-free. Δεν πρόκειται να ασχοληθούμε με την θεωρία τυπικών (formal) γλωσσών προγραμματισμού, μια και αυτό είναι το αντικείμενο της Θεματικής Υποενότητας ‘Αυτόματα και Τυπικές Γλώσσες’ αλλά θα αναφέρουμε μόνο εκείνα τα στοιχεία που είναι απολύτως απαραίτητα για να κατανοήσετε την λειτουργία της Συντακτικής Ανάλυσης.

## ΕΝΟΤΗΤΑ 3.1 ΓΡΑΜΜΑΤΙΚΕΣ CONTEXT-FREE

Είναι σύνηθες να ορίζουμε μερικές δομές γλωσσών προγραμματισμού με ένα αναδρομικό τρόπο (recursively). Για παράδειγμα, εάν  $E$  σημαίνει **έκφραση** (expression) και  $S_i$  σημαίνει **εντολή** (statement) τότε τα παρακάτω είναι επίσης εντολές (τα  $E_i$  είναι επίσης εκφράσεις).

**if**  $E$  **then**  $S_1$  **else**  $S_2$  (1)

**begin**  $S_1$  ;  $S_2$  ; ... ;  $S_n$  **end** (2)

Μια έκφραση είναι και η

$E_1 + E_2$  (3)

Αν χρησιμοποιήσουμε τις συντακτικές κατηγορίες "statement", και "expression" στην κλάση των εντολών και των εκφράσεων αντίστοιχα, τότε η (1) μπορεί να γραφτεί σύμφωνα με τον **Κανόνα** (Production -Rewriting Rule):

statement  $\rightarrow$  **if** expression **then** statement **else** statement (4)

η δε (3) μπορεί να γραφτεί σύμφωνα με τον κανόνα:

expression  $\rightarrow$  expression + expression

Η (2) όμως παρουσιάζει πρόβλημα αν γραφτεί σαν:

statement  $\rightarrow$  **begin** statement ; ... ; statement **end**

διότι χρειαζόμαστε κάθε κανόνα να έχει ακριβή αριθμό από σύμβολα.

Για να γραφτεί η (2) σωστά χρειάζεται να εισάγουμε άλλη μια συντακτική κατηγορία την οποία θα ονομάσουμε "statement-list". Η (2) τώρα μπορεί να περιγραφεί από τους εξής δυο κανόνες:

statement  $\rightarrow$  **begin** statement-list **end**

statement-list  $\rightarrow$  statement  
| statement; statement-list (5)

Όπου το σύμβολο "|" σημαίνει είτε. Δηλαδή ο δεύτερος κανόνας της (5) έχει δυο εναλλακτικά δεξιά μέρη.

Ένα σύνολο από κανόνες όπως η (5) είναι ένα παράδειγμα γραμματικής. Μια γραμματική περιλαμβάνει 4 οντότητες: **Τερματικά** σύμβολα, **Μη Τερματικά** σύμβολα, το **Αρχικό σύμβολο** και ένα σύνολο από **Κανόνες** (Terminals, Nonterminals, Start Symbol, Productions).

Τερματικά σύμβολα είναι τα βασικά σύμβολα που σχηματίζουν τις συμβολοσειρές της γλώσσας. Τα tokens και τα τερματικά σύμβολα είναι συνώνυμα στις γλώσσες προγραμματισμού. Π.χ. **begin**, **end**, **;**, **+** είναι τερματικά σύμβολα.

Μη Τερματικά σύμβολα (συνώνυμα με τις συντακτικές κατηγορίες και τις συντακτικές μεταβλητές), είναι ειδικά σύμβολα που υποδηλώνουν συμβολοσειρές οι οποίες εκφράζουν την συντακτική δομή της γλώσσας (πχ το statement στον κανόνα: statement  $\rightarrow$  **begin** statement-list **end**).

Ένα από τα Μη Τερματικά σύμβολα έχει επιλεγεί σαν Αρχικό σύμβολο και υποδηλώνει την γλώσσα.

Κανόνες είναι οι συντακτικοί κανόνες που καθορίζουν τους τρόπους με τους οποίους οι συντακτικές κατηγορίες μπορούν να κατασκευαστούν ή μια από την άλλη και από τα Τερματικά σύμβολα. Οι κανόνες συνίστανται από ένα Μη Τερματικό σύμβολο ακολουθούμενο από ένα  $\rightarrow$  (ή

::=) και ακολουθούμενο από μια συμβολοσειρά αποτελούμενη από Τερματικά και Μη Τερματικά σύμβολα.

Με τον συμβολισμό BNF (Backus-Naur-Form), η (5) περιγράφεται από τους παρακάτω κανόνες

statement  $\rightarrow$  **begin** statement-list **end**  
 statement-list  $\rightarrow$  statement  
 statement-list  $\rightarrow$  statement; statement-list

**Παράδειγμα 1 / Κεφ. 3.** Μια Context-free γραμματική η οποία μπορεί να περιγράψει απλές αριθμητικές εκφράσεις είναι η (6) η οποία έχει σαν Μη Τερματικά σύμβολα τα E και A, με το E να είναι το Αρχικό σύμβολο. Τα δε Τερματικά σύμβολα είναι τα:

**id** + - \* /  $\uparrow$  ( )

Τέλος οι Κανόνες είναι οι:

$E \rightarrow E A E \mid (E) \mid - E \mid \text{id}$   
 $A \rightarrow + \mid - \mid * \mid / \mid \uparrow$

(6)

Οι περισσότεροι συμβολισμοί που θα συναντήσουμε στη συνέχεια του βιβλίου είναι προφανείς, εκτός των εξής συμβάσεων:

- Κεφαλαία γράμματα από το τέλος της Λατινικής αλφαβήτου (X,Y,Z κλπ) παριστάνουν είτε Μη Τερματικά είτε Τερματικά.
- Μικρά γράμματα από το τέλος της Λατινικής αλφαβήτου (u,v,...,z) παριστάνουν συμβολοσειρές από Τερματικά.
- Μικρά Ελληνικά (α,β,γ) παριστάνουν συμβολοσειρές από σύμβολα της γραμματικής. (δηλαδή συμβολοσειρές από Τερματικά ή/και Μη Τερματικά σύμβολα) Π.χ.  $A \rightarrow a_1$

Η σειρά από τους Κανόνες  
 $A \rightarrow a_1 \mid A \rightarrow a_2, \dots, A \rightarrow a_k$  (A-Κανόνες)  
 συντομεύεται στην  
 $A \rightarrow a_1 \mid a_2 \mid \dots \mid a_k$ , όπως έγινε στο (6).

### 3.1.1 ΔΥΝΑΤΟΤΗΤΕΣ ΤΩΝ CONTEXT-FREE ΓΡΑΜΜΑΤΙΚΩΝ (CF)

Οι Context-free γραμματικές μπορούν να περιγράψουν την περισσότερη, αλλά όχι ολόκληρη την σύνταξη των γλωσσών προγραμματισμού (A. V. Aho, R. Sethi, J. D. Ullman, Compilers: Principles techniques, and Tools, Addison-Wesley, 1986). Είδαμε ότι οι Κανονικές Εκφράσεις περιγράφουν την σύνταξη των tokens. Οποιαδήποτε συντακτική δομή περιγράφεται από Κανονικές Εκφράσεις μπορεί επίσης να περιγραφεί και από κάποια Context-free γραμματική.

**Παράδειγμα 2 / Κεφ. 3.** Κανονική Έκφραση:

$(a \mid b)(a \mid b \mid 0 \mid 1)^*$

και η CF γραμματική:

$S \rightarrow aA \mid bA$   
 $A \rightarrow aA \mid bA \mid 0A \mid 1A \mid \varepsilon$

περιγράφουν την ίδια γλώσσα διότι και οι δυο δημιουργούν τις ίδιες συμβολοσειρές τερματικών συμβόλων.

Διάφοροι λόγοι μας αναγκάζουν να χρησιμοποιούμε Κανονικές Εκφράσεις αν και θα μπορούσαμε να χρησιμοποιήσουμε Context-free γραμματικές αντ' αυτών. Πρώτον οι λεκτικοί κανόνες είναι απλοί και δεν χρειάζεται συμβολισμός τύπου Context-free γραμματικών οι οποίοι είναι δυσκολότεροι να γίνουν αντιληπτοί. Δεύτερον είναι ευκολότερο να κατασκευάσουμε αναγνωριστές (recognizers) από τις Κανονικές Εκφράσεις παρά από Context-free γραμματικές.

Δεν υπάρχουν ακριβείς οδηγίες για το τι πρέπει να περιληφθεί στους λεκτικούς κανόνες και τι στους συντακτικούς. Οι Κανονικές Εκφράσεις είναι πιο κατάλληλες για να περιγράφουν τη δομή των λεκτικών κατασκευών (identifiers, constants, keywords, κλπ), ενώ οι Context-free γραμματικές είναι απαραίτητες για την περιγραφή απεριόριστα φωλιασμένων δομών (**begin - end**, **if - then - else**, κλπ).

### Άσκηση Αυτοαξιολόγησης 1 / Κεφ. 3

Μπορείτε να εξηγήσετε πού και γιατί χρησιμοποιούνται οι Κανονικές Εκφράσεις και πού οι Context Free γραμματικές στις γλώσσες προγραμματισμού;

## ΕΝΟΤΗΤΑ 3.2 ΔΗΜΙΟΥΡΓΙΑ ΤΩΝ ΔΕΝΔΡΩΝ ΑΝΙΧΝΕΥΣΗΣ

Στην γραμματική  $E \rightarrow E + E \mid E * E \mid (E) \mid -E \mid id$  (7)

όπου  $E$  σημαίνει έκφραση, ο Κανόνας  $E \rightarrow -E$  δηλώνει ότι μια έκφραση με μείον μπροστά της είναι επίσης μια έκφραση. Αυτό περιγράφεται από την σχέση  $E \Rightarrow -E$  που διαβάζεται: "η  $E$  παράγει (derives) -  $E$ ".

Η  $E \rightarrow (E)$  λέει ότι μπορούμε να αντικαταστήσουμε μια εμφάνιση της  $E$  σε οποιαδήποτε συμβολοσειρά από σύμβολα της γραμματικής, με  $(E)$ .

Μια άλλη περίπτωση είναι η:

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

μια τέτοια ακολουθία αντικαταστάσεων ονομάζεται **παραγωγή** (derivation) του **-(id)** από το  $E$ .

Σε πιο αφηρημένη μορφή λέμε ότι  $\alpha\beta \Rightarrow \alpha\gamma$  αν  $A \rightarrow \gamma$  είναι ένας Κανόνας και  $\alpha$  και  $\beta$  είναι αυθαίρετες συμβολοσειρές από σύμβολα της γραμματικής.

Αν  $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ , λέμε ότι η  $\alpha_1$  **παράγει**  $\alpha_n$ .

Το σύμβολο  $\Rightarrow$  μεταφράζεται σε "παράγει σε ένα βήμα", ενώ το  $\Rightarrow^*$  μεταφράζεται σε "παράγει σε μηδέν ή περισσότερα βήματα".

Έτσι ισχύει ότι:

1.  $\alpha \Rightarrow^* \alpha \quad \forall$  συμβολοσειρά  $\alpha \quad (\forall = \text{για κάθε})$

2. αν  $\alpha \Rightarrow^* \beta$  και  $\beta \Rightarrow \gamma$ , τότε  $\alpha \Rightarrow^* \gamma$

Ακόμη  $\alpha \Rightarrow^+ \beta$  σημαίνει "το  $\alpha$  παράγει  $\beta$  σε ένα ή περισσότερα βήματα".

Δοθείσης μιας context-free γραμματικής  $G$  με Αρχικό σύμβολο  $S$  μπορούμε να χρησιμοποιήσουμε την σχέση  $\Rightarrow^+$  για να ορίσουμε την  **$L(G)$  (την γλώσσα που δημιουργείται από την  $G$ )**. Οι συμβολοσειρές της  $L(G)$  μπορούν να περιέχουν μόνο Τερματικά σύμβολα της  $G$ . Λέμε ότι μια συμβολοσειρά  $w$  από Τερματικά ανήκει στην  $L(G)$  εάν και μόνο εάν  $S \Rightarrow^+ w$ . Η

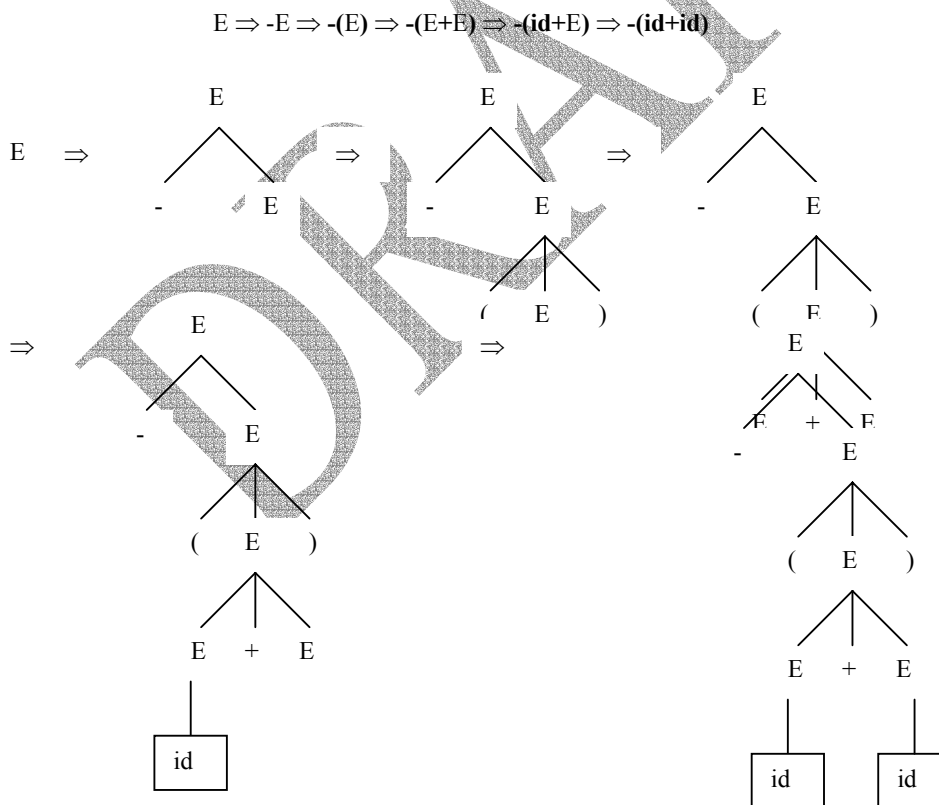
συμβολοσειρά  $w$  καλείται **πρόταση** (sentence) της  $G$ . Εάν  $S \Rightarrow^* a$ , όπου το  $a$  μπορεί να περιέχει και Μη Τερματικά, τότε λέμε ότι το  $a$  είναι μια **προτασιακή μορφή** (sentential form) της γραμματικής  $G$ .

### Άσκηση Αυτοαξιολόγησης 2 / Κεφ. 3

Μπορείτε να εξηγήσετε γιατί η συμβολοσειρά  $-(id+ id)$  είναι μια από τις προτάσεις της γλώσσας που προκύπτει από την γραμματική (7);

Για να καταλάβουμε καλύτερα πώς λειτουργούν μερικοί Συντακτικοί Αναλυτές θα θεωρήσουμε παραγωγές (derivations) στις οποίες μόνο το αριστερότερο Μη Τερματικό σύμβολο μιας προτασιακής μορφής αντικαθίσταται σε κάθε βήμα. Τέτοιες παραγωγές καλούνται **αριστερές παραγωγές** (Leftmost Derivations). Ανάλογα ορίζονται και **οι δεξιές παραγωγές** (Rightmost Derivations).

Μια γραφική απεικόνιση των παραγωγών είναι τα **Δένδρα Ανίχνευσης**. Σαν παράδειγμα ας δούμε την παραγωγή  $E \Rightarrow^+ -(id+ id)$  (8) με βήματα. Κατ' αρχήν θεωρούμε το  $E$  σαν ρίζα του αρχικού δένδρου και εφαρμόζοντας με σειρά τους κανόνες  $E \rightarrow -E$ ,  $E \rightarrow (E)$ ,  $E \rightarrow E+E$ ,  $E \rightarrow id$ , και  $E \rightarrow id$  παίρνουμε την παρακάτω σειρά από παραγωγές οι οποίες αντιστοιχούν στην τμηματική κατασκευή του δένδρου ανίχνευσης όπως δείχνεται στο Σχήμα 3.1.

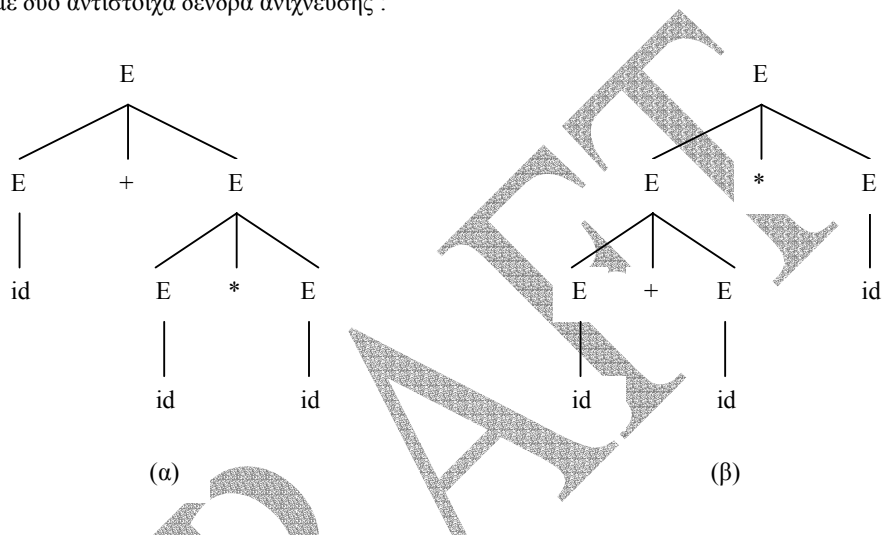


Σχήμα 3.1 Αριστερή παραγωγή του δένδρου ανίχνευσης της πρότασης  $-(id+id)$

**Δραστηριότητα 1 / Κεφ. 3** Η πρόταση  $id+id*id$  της γραμματικής (7), έχει δύο διαφορετικές αριστερές παραγωγές. Μπορείς να εξηγήσεις γιατί;

$  \begin{aligned}  E &\Rightarrow E + E \\  &\Rightarrow id + E \\  &\Rightarrow id + E * E \\  &\Rightarrow id + id * E \\  &\Rightarrow id + id * id  \end{aligned}  $	$  \begin{aligned}  E &\Rightarrow E * E \\  &\Rightarrow E + E * E \\  &\Rightarrow id + E * E \\  &\Rightarrow id + id * E \\  &\Rightarrow id + id * id  \end{aligned}  $
---	--

με δύο αντίστοιχα δένδρα ανίχνευσης :



Σχήμα 3.2 Τα δένδρα ανίχνευσης της πρότασης  $id+id*id$

Πράγματι αν ξεκινήσετε με τον κανόνα  $E \rightarrow E+E$  θα καταλήξετε στο δένδρο (α) που αντιστοιχεί στην άνω αριστερά παραγωγή του Σχήματος 3.2 ενώ αν ξεκινήσετε με τον κανόνα  $E \rightarrow E*E$  θα καταλήξετε στο δένδρο (β) το οποίο αντιστοιχεί στην άνω δεξιά παραγωγή του ίδιου σχήματος. Ο λόγος για τον οποίο η πρόταση  $id+id*id$  έχει δύο διαφορετικά δένδρα ανίχνευσης είναι διότι η γραμματική (7) είναι **διφορούμενη**. Αν παρατηρήσετε λίγο την γραμματική αυτή θα δείτε ότι δεν εκφράζει τις προτεραιότητες ούτε τις προσεταιριστικότητες των τελεστών. Για παράδειγμα η γραμματική αυτή υποδηλώνει ότι οι τελεστές + και \* έχουν την ίδια προτεραιότητα. Επιπλέον δεν έχουμε τρόπο να αποφασίσουμε ποιόν από τους εναλλακτικούς κανόνες πρέπει να χρησιμοποιήσουμε σε κάθε παραγωγή. Έτσι μπορούμε να φτιάξουμε είτε το ένα δένδρο είτε το άλλο.

### ΕΝΟΤΗΤΑ 3.3 ΔΙΦΟΡΟΥΜΕΝΕΣ ΓΡΑΜΜΑΤΙΚΕΣ

**Διφορούμενη/ασαφής** λέγεται μια γραμματική η οποία παράγει περισσότερα από ένα δένδρα ανίχνευσης, για μια πρόταση. Για ορισμένους τύπους συντακτικών αναλυτών, θέλουμε η γραμματική να γίνει σαφής γιατί αλλιώς δεν μπορούμε να προσδιορίσουμε ένα μοναδικό δένδρο για μια πρόταση. Σε μερικές περιπτώσεις γίνεται χρήση Διφορούμενων γραμματικών μαζί με κανόνες αποσαφήνισης.

Έστω η γραμματική:  $E \rightarrow E+E \mid E-E \mid E * E \mid E / E \mid E \uparrow E \mid (E) \mid -E \mid \text{id}$  (9)

Η γραμματική αυτή όπως και η (7) είναι διαφορούμενη. Μπορούμε όμως να απο-διφοροποιήσουμε και τις δύο, καθορίζοντας προσηταιριστικότητα και προτεραιότητα στους τελεστές. Ας υποθέσουμε ότι θέλουμε να τους δώσουμε τις γνωστές από τα Μαθηματικά προτεραιότητες:

- (unary minus)
↑
* /
+ -

και επίσης να ορίσουμε τον τελεστή  $\uparrow$  να είναι δεξιά προσηταιριστικός [δηλαδή,  $a \uparrow b \uparrow c = a \uparrow (b \uparrow c)$ ], και τους υπόλοιπους τελεστές να είναι αριστερά-προσηταιριστικοί. [δηλαδή,  $a-b-c = (a-b)-c$ ].

Σύμφωνα με τους κανόνες αυτούς το δένδρο ( $\beta$ ) του Σχήματος 3.2 δεν είναι σωστό για την πρόταση  $\text{id}+\text{id}*\text{id}$ , μια και εδώ εμφανίζεται το  $+$  να έχει μεγαλύτερη προτεραιότητα από το  $*$ .

Μπορούμε να ξαναγράψουμε τις γραμματικές ενσωματώνοντας σε αυτές τους κανόνες προσηταιρισμού και προτεραιότητας.

### Παράδειγμα 3 / Κεφ. 3

Ας πάρουμε την διαφορούμενη γραμματική (9). Κατ' αρχήν εισάγουμε ένα Μη Τερματικό σύμβολο για κάθε επίπεδο προτεραιότητας.

Ονομάζουμε *Element* μια υποέκφραση η οποία είναι μη υποδιαίρεσιμη άρα είναι είτε ένας identifier ή μια έκφραση με παρενθέσεις. Έχουμε λοιπόν τον κανόνα:

$$\text{element} \rightarrow (\text{expression}) \mid \text{id}$$

Κατόπιν εισάγουμε την κατηγορία των *primaries* τα οποία είναι elements με ένα ή περισσότερους από τους τελεστές με την μεγαλύτερη προτεραιότητα:

$$\text{primary} \rightarrow \text{primary} \mid \text{element}$$

Κατόπιν εισάγουμε τα *factors* σαν ακολουθίες από ένα ή περισσότερα primaries συνδεδεμένα με τον τελεστή ύψωσης σε δύναμη ( $\uparrow$ )

$$\text{factor} \rightarrow \text{primary} \uparrow \text{factor} \mid \text{primary}$$

Σημειώνουμε ότι η εκλογή  $\text{primary} \uparrow \text{factor}$ , αντί του  $\text{factor} \uparrow \text{primary}$  κάνει τις εκφράσεις  $a \uparrow b \uparrow c$  να ομαδοποιούνται σαν  $a \uparrow (b \uparrow c)$ .

Στη συνέχεια εισάγουμε το *term*, που είναι ακολουθία από ένα ή περισσότερα factors που συνδέονται με τους τελεστές πολλαπλασιασμού ( $*$ ,  $/$ ).

$$\text{term} \rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor}$$

Οι κανόνες αυτοί αναγκάζουν τα terms να ομαδοποιούνται από τα αριστερά [δηλαδή,  $a * b * c = (a * b) * c$ ].

Τέλος τα *expressions* είναι ακολουθίες από ένα ή περισσότερα terms και συνδέονται με τους προσθετικούς τελεστές ( $+$ ,  $-$ ).

$$\text{Expression} \rightarrow \text{expression} + \text{term} \mid \text{expression} - \text{term} \mid \text{term}$$

Και εδώ ο κανόνας επιβάλλει αριστερή προσηταιριστικότητα.

Έτσι η απο-διφοροποιημένη γραμματική είναι η:

$$\begin{aligned}
 \text{expression} &\rightarrow \text{expression} + \text{term} & (10) \\
 &| \text{expression} - \text{term} \\
 &| \text{term} \\
 \text{term} &\rightarrow \text{term} * \text{factor} \\
 &| \text{term} / \text{factor} \\
 &| \text{factor} \\
 \text{factor} &\rightarrow \text{Primary} \uparrow \text{factor} \\
 &| \text{Primary} \\
 \text{Primary} &\rightarrow -\text{Primary} \\
 &| \text{element} \\
 \text{element} &\rightarrow (\text{expression}) \\
 &| \text{id}
 \end{aligned}$$

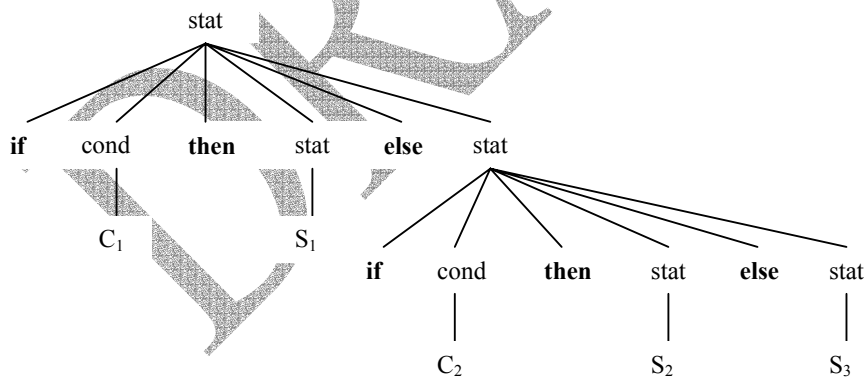
**Δραστηριότητα 2 / Κεφ. 3** Το παρακάτω κομμάτι μιάς Context Free γραμματικής δημιουργεί εντολές υπό συνθήκη. Όμως η γραμματική αυτή είναι διφορούμενη. Μπορείτε να εξηγήσετε γιατί και στη συνέχεια να την αποδιοφοροποιήσετε; Το πρώτο μέρος είναι εύκολο, το δεύτερο όμως είναι μάλλον δύσκολο γιατί να μην απογοητευθείτε αν δεν τα καταφέρετε.

$$\begin{aligned}
 \text{stat} &\rightarrow \text{if cond then stat} \\
 &| \text{if cond then stat else stat} & (11) \\
 &| \text{other - stat}
 \end{aligned}$$

Για να δούμε αν η γραμματική (11) είναι διφορούμενη πρέπει να εξετάσουμε μήπως για κάποια πρότασή της προκύπτουν περισσότερα από ένα δένδρα αντίχνευσης. Έτσι βλέπουμε ότι η πρόταση

**if C<sub>1</sub> then S<sub>1</sub> else if C<sub>2</sub> then S<sub>2</sub> else S<sub>3</sub>**

θα δώσει το παρακάτω δένδρο:



Σχήμα 3.3 Δένδρο αντίχνευσης της **if C<sub>1</sub> then S<sub>1</sub> else if C<sub>2</sub> then S<sub>2</sub> else S<sub>3</sub>**

Η (11) όμως είναι τελικά διφορούμενη διότι η πρόταση

$$\text{if C}_1 \text{ then if C}_2 \text{ then S}_1 \text{ else S}_2 \quad (12)$$

έχει τα δύο δένδρα αντίχνευσης που δίνονται στο Σχήμα 3.4.

Σε όλες τις γλώσσες προγραμματισμού που περιέχουν τέτοιες εντολές προτιμάται το δένδρο (α). Ο γενικός κανόνας που χρησιμοποιείται είναι: κάθε **else** να ζευγαρώνεται με το πλησιέστερο προηγούμενο αζευγάρωτο **then**. Αυτό αποτελεί τον κανόνα αποδιοφοροποίησης ο οποίος, αν θέλουμε, μπορεί να εισαχθεί στην γραμματική (11) και να πάρουμε την εξής:



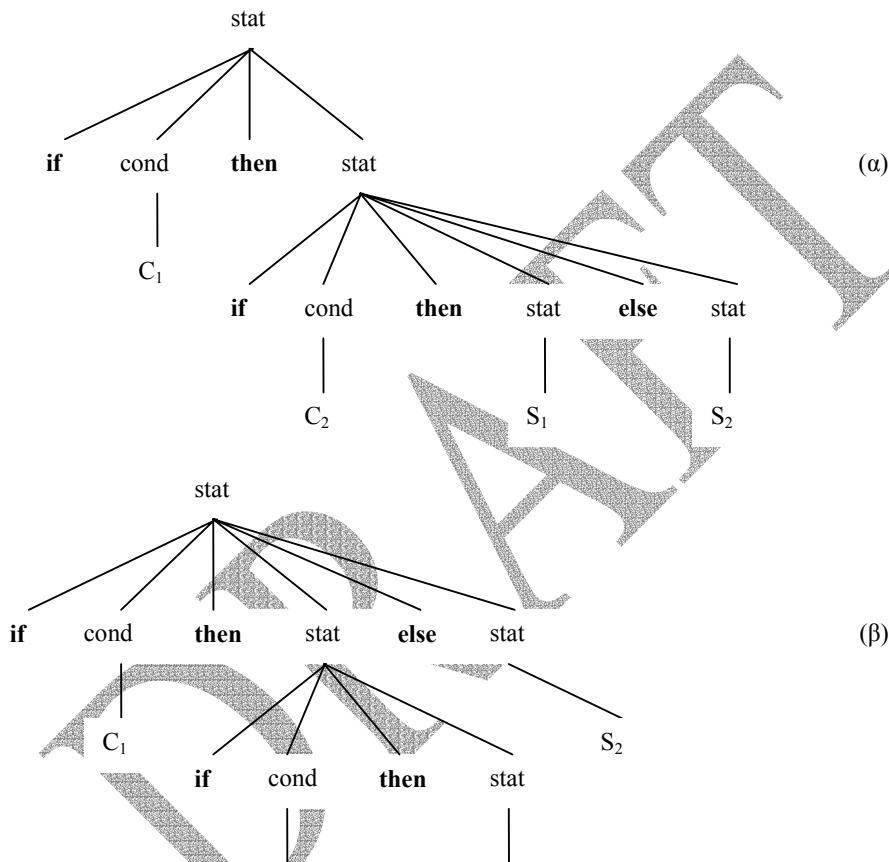
$stat \rightarrow matched - stat \mid unmatched - stat$

$matched - stat \rightarrow \mathbf{if} \text{ cond } \mathbf{then} \text{ matched - stat } \mathbf{else} \text{ matched - stat}$   
 $\mid \text{ other - stat}$

$unmatched - stat \rightarrow \mathbf{if} \text{ cond } \mathbf{then} \text{ stat}$

$\mid \mathbf{if} \text{ cond } \mathbf{then} \text{ matched - stat } \mathbf{else} \text{ unmatched - stat}$

Η γραμματική αυτή δημιουργεί τις ίδιες συμβολοσειρές με την (11) αλλά επιτρέπει ένα μόνο δένδρο για την (12), δηλαδή αυτό που συσχετίζει κάθε **else** με το προηγούμενο αζευγάρωτο **then**. Στο Κεφάλαιο 4 θα δούμε και ένα άλλο τρόπο αποδιφοροποίησης της (11).



Σχήμα 3.4 Τα δένδρα ανίχνευσης της  $\mathbf{if} \ C_1 \ \mathbf{then} \ \mathbf{if} \ C_2 \ \mathbf{then} \ S_1 \ \mathbf{else} \ S_2$

### Άσκηση Αυτοαξιολόγησης 3 / Κεφ. 3

Σας δίνεται η διαφορούμενη γραμματική ( όπου **con** υποδηλώνει constant )

$E \rightarrow E+E \mid E * E \mid (E) \mid \mathbf{id} \mid \mathbf{con}$

Προσπαθείστε να την αποδιφοροποιήσετε εφαρμόζοντας την μεθοδολογία που περιγράψαμε στην ενότητα 3.3.

## ΕΝΟΤΗΤΑ 3.4 ΙΕΡΑΡΧΙΑ ΓΛΩΣΣΩΝ ΚΑΤΑ CHOMSKY

Ο Chomsky καθόρισε τέσσερις βασικές κλάσεις γλωσσών μέσω γραμματικών, οι οποίες είναι τετράδες (V,T,P,Z) όπου:

1. V είναι η αλφάβητος (Τερματικά και Μη Τερματικά σύμβολα)

2.  $T, (T \subseteq V)$  είναι το σύνολο των τερματικών συμβόλων
3.  $P$  είναι ένα πεπερασμένο σύνολο από κανόνες (επαναγραφής) και
4.  $Z$  είναι το αρχικό ή διακεκριμένο σύμβολο της γραμματικής.

Η γλώσσα μιας γραμματικής είναι το σύνολο των τερματικών ακολουθιών (terminal strings) τα οποία μπορούν να δημιουργηθούν από το αρχικό σύμβολο  $Z$ .

Η διαφορά μεταξύ των τεσσάρων τύπων γραμματικών ευρίσκεται στην μορφή των κανόνων επαναγραφής που επιτρέπονται στο  $P$ .

Μία γραμματική  $G$  είναι (Chomsky) **τύπου 0** ή **phrase structure** γραμματική αν οι κανόνες του  $P$  έχουν την μορφή

$$(0) u \rightarrow v, \quad \text{με } u \in V^+ \text{ και } v \in V^*$$

Δηλαδή, το αριστερό  $u$  μπορεί να είναι μια ακολουθία συμβόλων και το δεξιό μέρος να είναι άδειο.

Μια γραμματική  $G$  είναι **τύπου 1** ή **context sensitive** (εξαρτώμενη από το περιβάλλον κείμενο) όταν περιορίσουμε τους κανόνες επαναγραφής στην μορφή:

$$(1) xUy \rightarrow xuy, \quad \text{με } U \in V-T, x, y \in V^*, \text{ και } u \in V^+$$

Ο όρος context sensitive αναφέρεται στο γεγονός ότι επιτρέπεται να ξαναγράψουμε το  $U$  σαν  $u$  μόνο μέσα στο περιβάλλον  $x...y$ .

Μια γραμματική  $G$  είναι **τύπου 2** ή **Context-Free** (ανεξάρτητη περιβάλλοντος) αν όλοι οι κανόνες έχουν την μορφή:

$$(2) U \rightarrow u, \quad \text{με } U \in V-T \text{ και } u \in V^*$$

Η κλάση αυτή καλείται ανεξάρτητη περιβάλλοντος διότι το  $U$  μπορεί να ξαναγραφεί σαν  $u$  ανεξάρτητα από το υπόλοιπο περιβάλλον κείμενο στο οποίο ευρίσκεται το  $U$ . Όπως φαίνεται από τον ορισμό σε μια CF γραμματική ένας κανόνας μπορεί να έχει την μορφή  $U \rightarrow \epsilon$  ( $\epsilon$ =κενό). Δοθείσης μιας CF γραμματικής  $G$ , μπορεί να κατασκευασθεί μια  **$\epsilon$ -free** γραμματική  $G1$  (η οποία δεν περιέχει κανόνα της μορφής  $U \rightarrow \epsilon$ ), τέτοια ώστε  $L(G1) = L(G) - \{\epsilon\}$ . Επιπλέον αν η  $G$  είναι μη διαφορούμενη τέτοια είναι και η  $G1$ .

Περιορίζοντας τους κανόνες ακόμη περισσότερο, στην μορφή :

$$(3) U \rightarrow N \text{ ή } U \rightarrow WN, \quad \text{με } N \in T, \text{ και } U \text{ και } W \in V-T$$

έχουμε μια γραμματική **τύπου 3** ή Κανονική γραμματική (Regular grammar).

## Παράδειγμα 4 / Κεφ. 3 Παραδείγματα κανόνων από Κανονικές Γραμματικές

Παρατηρήστε ότι τα Μη Τερματικά σύμβολα μπορούν να περιβάλλονται από “<” και “>” για λόγους καλύτερης αναγνωσιμότητας των κανόνων.

<s-delim>  $\rightarrow + | - | ( | )$

<d-delim>  $\rightarrow </> | </>^* | <ast>^* | <colon> =$

<slash>  $\rightarrow /$

<ast>  $\rightarrow *$

<colon>  $\rightarrow :$

## Άσκηση Αυτοαξιολόγησης 4 / Κεφ. 3

Οι Κανονικές Εκφράσεις για ονόματα μεταβλητών και για ακέραιες σταθερές είναι  $l(l|d)^*$  και  $d(d)^*$  αντίστοιχα (  $l$  υποδηλώνει letter και  $d$  υποδηλώνει digit). Προσπαθήστε να τις γράψετε με την βοήθεια των Κανονικών γραμματικών.

## ΣΥΝΟΨΗ ΚΕΦΑΛΑΙΟΥ

Στο κεφάλαιο αυτό είδαμε μια άτυπη περιγραφή των γραμματικών τύπου Context Free που είναι οι γραμματικές οι οποίες κατά κύριο λόγο χρησιμοποιούνται για τον συντακτικό ορισμό των γλωσσών προγραμματισμού. Αναφέραμε βέβαια και τους τυπικούς ορισμούς του Chomsky για λόγους πληρότητας αλλά και για να αναφερθούμε και στις Κανονικές γραμματικές με τις οποίες είδαμε ότι μπορούμε να εκφράσουμε τις λεκτικές δομές των γλωσσών προγραμματισμού. Επειδή οι Κανονικές γραμματικές είναι στην ουσία υποσύνολο των Context Free γραμματικών θα μπορούσαμε να πούμε ότι με τις Context Free γραμματικές μπορούμε να εκφράσουμε ολόκληρη την συντακτική και λεκτική δομή των γλωσσών.

Δώσαμε έμφαση στις παραγωγές από το Αρχικό σύμβολο της γραμματικής και ορίσαμε την πρόταση και την προτασιακή μορφή τις οποίες θα χρησιμοποιήσουμε στο επόμενο κεφάλαιο. Έμφαση, και μάλιστα ιδιαίτερη, δώσαμε στις διαφορούμενες γραμματικές και στο πώς τις αποδιφοροποιούμε.

## Απάντηση άσκησης 1 /κεφ. 3

Οι Κανονικές Εκφράσεις χρησιμοποιούνται στις γλώσσες προγραμματισμού για να περιγράψουν τις λεκτικές τους δομές, δηλαδή τα tokens. Αντίθετα, για να περιγράψουμε τις συντακτικές δομές της γλώσσας (εντολές) χρησιμοποιούμε τις Context Free γραμματικές. Μπορούμε φυσικά να χρησιμοποιήσουμε Context Free γραμματικές για να περιγράψουμε και τις λεκτικές δομές της γλώσσας, όμως επειδή οι Context Free γραμματικές είναι πολύ πιο ισχυρό εργαλείο από τις Κανονικές Εκφράσεις, είναι και πολύ πιο δύσκολο να φτιάξουμε ‘αναγνωριστές’ για Context Free γραμματικές απότι είναι για Κανονικές Εκφράσεις.

## Απάντηση άσκησης 2 /κεφ. 3

Μια συμβολοσειρά από τερματικά σύμβολα μιας γραμματικής είναι πρόταση της γλώσσας που δημιουργείται από την γραμματική αυτή εάν προκύπτει μέσα από μια σειρά παραγωγών οι οποίες ξεκινούν από το Αρχικό σύμβολο της γραμματικής. Έτσι, δεδομένου ότι το Αρχικό σύμβολο της γραμματικής (7) είναι το E, με την παρακάτω σειρά παραγωγών καταλήγουμε στην πρόταση  $-(id+id)$ .

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

Η παραγωγή αυτή είναι *δεξιότερη* παραγωγή. Στην ίδια πρόταση μπορείτε να καταλήξετε και με την *εξής αριστερότερη* παραγωγή

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

Οι συμβολοσειρές  $-E$ ,  $-(E)$ ,  $-(E+E)$ ,  $-(E+id)$  και  $-(id+E)$  είναι *προτασιακές* μορφές.

### Απάντηση άσκησης 3 /κεφ. 3

Ο λόγος που η γραμματική αυτή είναι διαφορούμενη είναι διότι έχοντας ένα μόνο κανόνα με πέντε εναλλακτικά δεξιά μέρη ορίζει ένα μόνο επίπεδο προτεραιότητας και έτσι οι τελεστές  $+$ ,  $*$  αλλά και οι μη υποδιαίρεσιμες εκφράσεις  $(E)$ , **id** και **con** έχουν όλα την ίδια προτεραιότητα. Θυμηθείτε ότι για να ορίσουμε τις προτεραιότητες εισάγουμε καινούρια Μη Τερματικά σύμβολα (τα οποία τα είχαμε αναφέρει ως *συντακτικές κατηγορίες*). Επίσης η γραμματική αυτή δεν εκφράζει ούτε την προσεταιριστικότητα των τελεστών. Και πάλι θυμηθείτε ότι την προσεταιριστικότητα την εκφράζουμε με τον τρόπο που γράφουμε τους κανόνες. Για παράδειγμα ο κανόνας

$$X \rightarrow X@Y$$

αποδίδει αριστερή προσεταιριστικότητα για τον τελεστή  $@$ , ενώ αντίθετα ο κανόνας

$$X \rightarrow Y@X$$

αποδίδει δεξιά προσεταιριστικότητα για τον τελεστή  $@$ .

Μετά από αυτά μπορούμε να προχωρήσουμε στην αποδιφοροποίηση.

Εισάγουμε πρώτα ένα νέο Μη Τερματικό σύμβολο **Elem** για τις μη διαίρεσιμες εκφράσεις με την υψηλότερη προτεραιότητα,

$$\text{Elem} \rightarrow (E) | \text{id} | \text{con}$$

Κατόπιν εισάγουμε το σύμβολο **Fact** για τον τελεστή πολλαπλασιασμού ο οποίος είναι αριστερά προσεταιριστικός και έχει την αμέσως χαμηλότερη προτεραιότητα.

$$\text{Fact} \rightarrow \text{Fact} * \text{Elem} | \text{Elem}$$

Τέλος, έχουμε το Αρχικό σύμβολο της γραμματικής  $E$  για τον τελεστή πρόσθεσης που έχει την χαμηλότερη προτεραιότητα και είναι επίσης αριστερά προσεταιριστικός.

$$E \rightarrow E + \text{Fact} | \text{Fact}$$

Έτσι η αρχική διαφορούμενη γραμματική

$$E \rightarrow E+E | E * E | (E) | \text{id} | \text{con}$$

Μετασχηματίστηκε στην αποδιφοροποιημένη μορφή

$$E \rightarrow E + \text{Fact} | \text{Fact}$$

$$\text{Fact} \rightarrow \text{Fact} * \text{Elem} | \text{Elem}$$

$$\text{Elem} \rightarrow (E) | \text{id} | \text{con}$$

### Απάντηση άσκησης 4 /κεφ. 3

Οι Κανονικές γραμματικές σύμφωνα με τον ορισμό τους έχουν την μορφή

$$U \rightarrow N \text{ ή } U \rightarrow WN, \quad \text{με } N \in T \text{ και } U \text{ και } W \in V-T$$

Έτσι, αν συμβολίσουμε το  $l$  με letter και το  $d$  με digit, για πιο ευανάγνωστους κανόνες, και επιπλέον χρησιμοποιήσουμε τα Μη Τερματικά σύμβολα  $\langle id \rangle$  και  $\langle int \rangle$  για να εκφράσουμε τα

ονόματα μεταβλητών και τις ακεραίες σταθερές αντίστοιχα τότε οι παρακάτω κανόνες είναι οι ζητούμενοι μια και παράγουν τις ίδιες συμβολοσειρές τερματικών συμβόλων.

`<id>`        → letter  
                 |`<id>` letter  
                 |`<id>`digit  
`<int>`        → digit  
                 |`<int>` digit

Παρατηρήστε ότι η Κανονική Έκφραση  $d(d)^*$  θα μπορούσε να είχε δοθεί και με τον ισοδύναμο ορισμό της  $d^+$  χωρίς φυσικά να αλλάξει τίποτα.

DRAFT