

# Análise sintática descendente preditiva. Parte I – análise recursiva

- Professor Eraldo Pereira Marinho
- Compiladores 2º semestre 2020
- UNESP/IGCE Rio Claro

# Noções preliminares – árvore de derivação

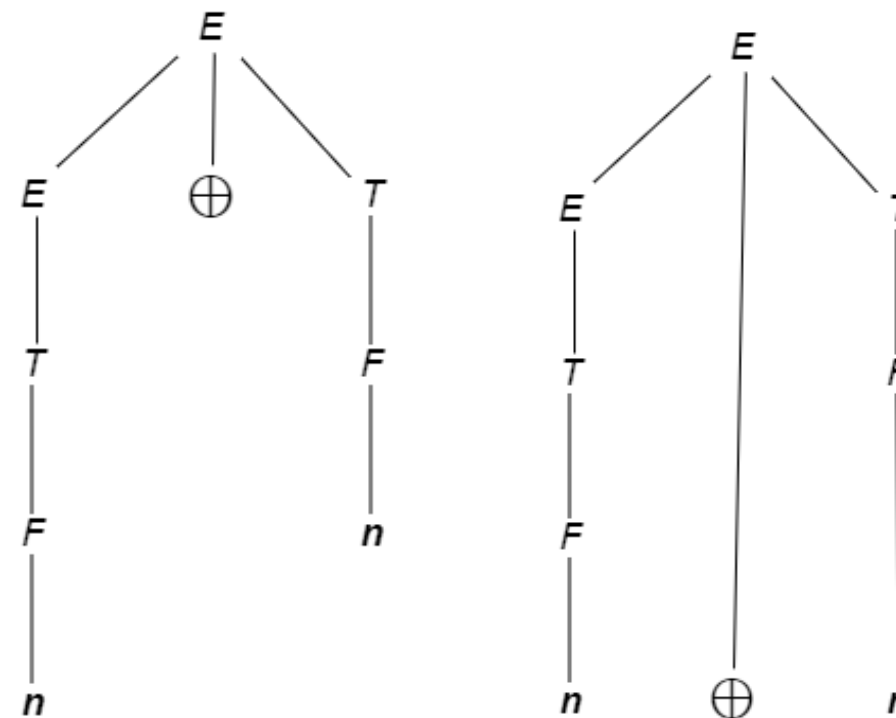
$$\begin{aligned} E &\rightarrow E \oplus T \mid T \\ T &\rightarrow T \otimes F \mid F \\ F &\rightarrow (E) \mid \mathbf{n} \end{aligned}$$

Derivação esquerda (usual):

$$\begin{aligned} E &\Rightarrow^e E \oplus T \Rightarrow^e T \oplus T \\ &\Rightarrow^e F \oplus T \Rightarrow^e \mathbf{n} \oplus T \\ &\Rightarrow^e \mathbf{n} \oplus F \Rightarrow^e \mathbf{n} \oplus \mathbf{n} \end{aligned}$$

Derivação direita (canônica):

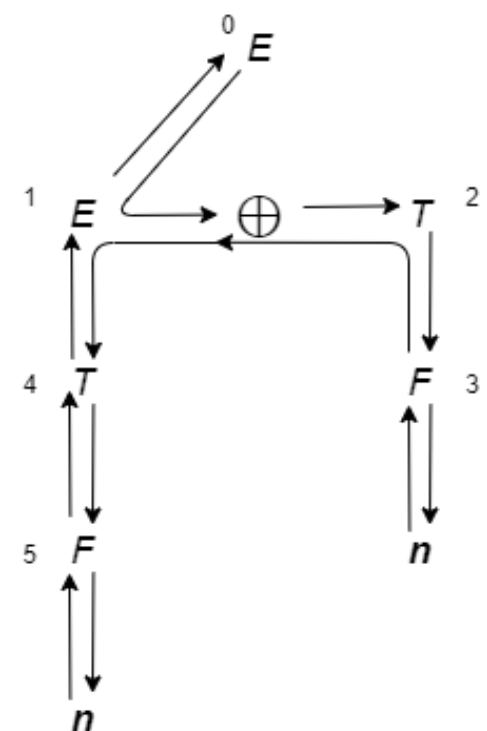
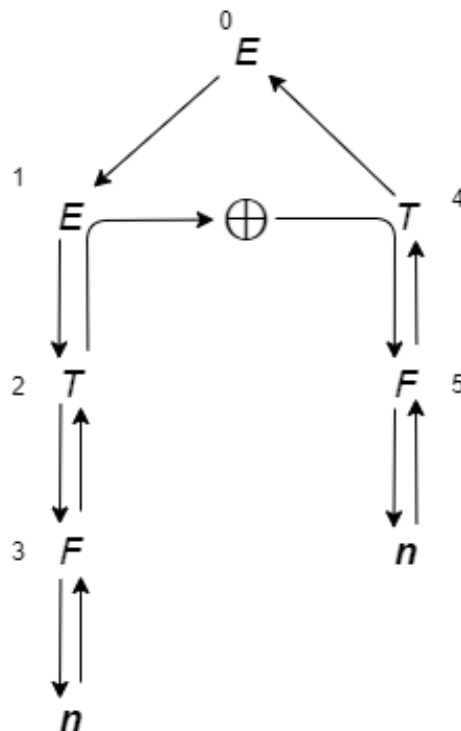
$$\begin{aligned} E &\Rightarrow^d E \oplus T \Rightarrow^d E \oplus F \\ &\Rightarrow^d E \oplus \mathbf{n} \Rightarrow^d T \oplus \mathbf{n} \\ &\Rightarrow^d F \oplus \mathbf{n} \Rightarrow^d \mathbf{n} \oplus \mathbf{n} \end{aligned}$$



# Noções preliminares – caminhamento na árvore de derivação

$$E^0 \Rightarrow E^1 \oplus T \Rightarrow T^2 \oplus T \Rightarrow F^3 \oplus T \\ \Rightarrow \mathbf{n} \oplus T^4 \Rightarrow \mathbf{n} \oplus F^5 \Rightarrow \mathbf{n} \oplus \mathbf{n}$$

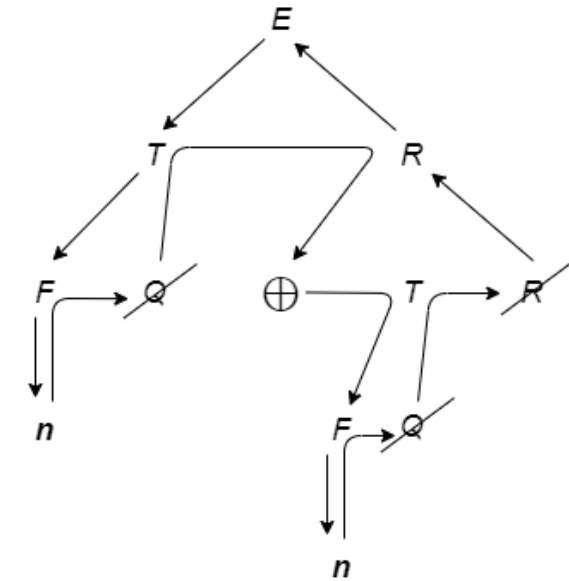
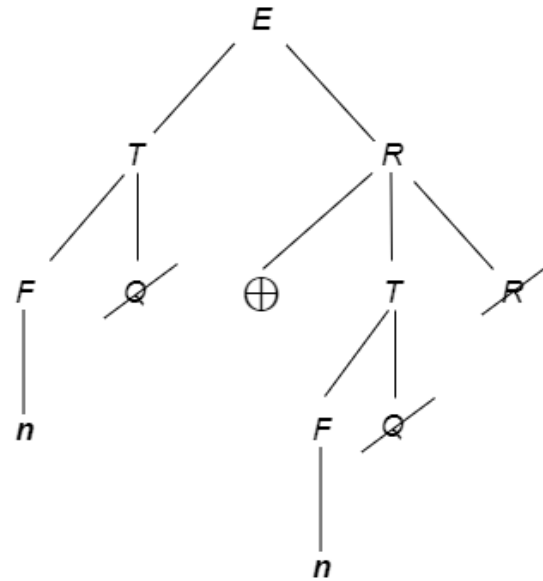
$$E^0 \Rightarrow^d E^1 \oplus T^2 \Rightarrow^d E^1 \oplus F^3 \\ \Rightarrow^d E^1 \oplus \mathbf{n} \Rightarrow^d T^4 \oplus \mathbf{n} \\ \Rightarrow^d F^5 \oplus \mathbf{n} \Rightarrow^d \mathbf{n} \oplus \mathbf{n}$$



# Noções preliminares – caminhamento na árvore de derivação

$$\begin{aligned}
 E &\rightarrow TR \\
 T &\rightarrow FQ \\
 R &\rightarrow \oplus TR \mid \varepsilon \\
 Q &\rightarrow \otimes FQ \mid \varepsilon \\
 F &\rightarrow (E) \mid n
 \end{aligned}$$

$$\begin{aligned}
 E &\Rightarrow TR \Rightarrow FQR \Rightarrow nQR \Rightarrow nR \\
 &\Rightarrow n \oplus TR \oplus FQR \Rightarrow n \oplus nQR \\
 &\Rightarrow n \oplus nR \Rightarrow n \oplus n
 \end{aligned}$$



# Noções preliminares – determinismo na escolha do lado direito

Percebemos que a gramática

$$\begin{aligned}E &\rightarrow TR \\T &\rightarrow FQ \\R &\rightarrow \oplus TR \mid \varepsilon \\Q &\rightarrow \otimes FQ \mid \varepsilon \\F &\rightarrow (E) \mid \mathbf{n}\end{aligned}$$

tem uma certa vantagem, se pensarmos em desenvolver uma máquina que emule o processo de derivação, em comparação com

$$\begin{aligned}E &\rightarrow E \oplus T \mid T \\T &\rightarrow T \otimes F \mid F \\F &\rightarrow (E) \mid \mathbf{n}\end{aligned}$$

Por que? A resposta é que, nesta última, uma tal máquina não tem como decidir de forma simples o que será escolhido, já da primeira produção.  $E \rightarrow E \oplus T$  ou  $E \rightarrow T$ ?

Suponhamos, novamente, que queremos derivar  $\mathbf{n} \oplus \mathbf{n}$  usando esta última gramática. Qual produção a máquina escolheria primeiramente se a única informação que essa teria era que o token lido antecipadamente seria  $\tau = \mathbf{n}$ ? Ou seja,  $E \Rightarrow T$  ou  $E \Rightarrow E \oplus T$ ? A máquina não saberia qual dessas produções escolheria. Daí viria uma alternativa que seria a heurística:  $E \Rightarrow T \Rightarrow F \Rightarrow \mathbf{n}$ , mas aí a máquina receberia um novo token, que seria  $\tau = \oplus$ . Aí, a máquina teria que devolver esse token e voltar a  $\tau = \mathbf{n}$  e escolher a segunda opção,  $E \Rightarrow E \oplus T \Rightarrow T \oplus T \Rightarrow F \oplus T \Rightarrow \mathbf{n} \oplus T \Rightarrow \mathbf{n} \oplus F \Rightarrow \mathbf{n} \oplus \mathbf{n}$ . O mesmo já não ocorreria com a primeira gramática.

# Forma normal de Greibach

- Os primeiros passos para um estudo sistemático da análise sintática é a análise sintática descendente recursiva, que requer uma forma normal de gramática livre do contexto, as gramáticas LL(1)
- Começamos pela definição da forma normal de Greibach:
- Uma gramática  $G = (N, \Sigma, S, P)$ , livre do contexto, é dita na forma normal de Greibach se e somente se suas produções são todas da forma  $A \rightarrow \sigma \alpha$ , onde  $\sigma \in \Sigma$  e  $\alpha \in (\Sigma \cup N)^*$ . Em outras palavras, o lado direito das produções deve ser explicitamente iniciados por terminais,  $\in \Sigma$
- O teorema de Greibach garante que toda boa gramática pode ser rescrita na forma normal de Greibach

# Teorema de Greibach

**Teorema 1.** Seja  $G = (N, \Sigma, S, P)$  uma boa gramática. Neste caso, é sempre possível apresentar uma gramática  $G' = (N', \Sigma, S, P')$ ,  $N \subseteq N'$  com suas produções  $P'$  na forma normal de Greibach.

**Prova.** Sendo  $G$  uma boa gramática, então não há produções nulas, nem produções unitárias e muito menos produções quebradas, sem vínculos. Contudo, pode haver produções na forma

$$A \rightarrow A\alpha \mid \beta; \quad A \in N, \alpha, \beta \in (\Sigma \cup N)^*, \beta \neq A\gamma$$

Por indução, temos que  $A$  pode derivar as seguintes situações:  $A \Rightarrow \beta$ ,  $A \Rightarrow A\alpha \Rightarrow \beta\alpha$  e, para um número  $n + 1$  de passos de derivação,  $A \Rightarrow^n A\alpha^n \Rightarrow \beta\alpha^n$ . Assim, podemos dizer que existe um  $R \in N'$  tal que  $R \Rightarrow^n \alpha^n$ , o que é possível se  $R \rightarrow \alpha R \mid \alpha$ . Neste caso, podemos substituir as produções acima pelas novas produções

$$A \rightarrow \beta R \mid \beta, \quad R \rightarrow \alpha R \mid \alpha$$

A demonstração ainda não está completa pois pode haver mais de uma produção com recursão esquerda em  $A$ . Suponhamos que existam as produções

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$$

É fácil generalizar o raciocínio anterior para mostrarmos que

$$A \Rightarrow^* A\alpha_{j_k} \dots \alpha_{j_2} \alpha_{j_1} \Rightarrow \beta_i \alpha_{j_k} \dots \alpha_{j_2} \alpha_{j_1}$$

Observemos que  $j_k$  é um inteiro qualquer em  $\{1, 2, \dots, m\}$  invocado no último nível de recursão (substituição), nível  $k$ . Desta última indução, temos que um novo conjunto de substituições será necessário para eliminarmos a recursão esquerda:

$$A \rightarrow \beta_1 R \mid \dots \mid \beta_n R \mid \beta_1 \mid \dots \mid \beta_n, \quad R \rightarrow \alpha_1 R \mid \dots \mid \alpha_m R \mid \alpha_1 \mid \dots \mid \alpha_m$$

Caso haja nova recursão esquerda nos  $\beta$ , o método acima é reaplicado até não mais recursão esquerda. Quando não há recursão esquerda, é possível obter, eventualmente,  $\beta_i$  na forma  $\beta_i = B\gamma_i$ .

# Teorema de Greibach - continuação

Dando continuidade à demonstração, produções do tipo  $A \rightarrow B\alpha$ , com  $B \neq^* A\gamma$ , devem ser recursivamente substituídas até não mais haver produções com lado direito iniciando com não terminal. Uma vez que não temos mais recursão esquerda, usamos o seguinte algoritmo:

**Entrada:** uma gramática LC sem recursão esquerda,  $G = (N, \Sigma, S, P)$

**Saída:** uma gramática na FNG,  $G' = (N, \Sigma, S, P')$

**Método:**

Inicialize  $P'$  com todas as produções  $P$  que já estejam na FNG

Repita

    Para toda produção de  $P$ , na forma  $A \rightarrow B\alpha$ , faça:

        Para toda produção  $B \rightarrow \beta$  de  $P'$ , faça:

            Acrescente  $A \rightarrow \beta\alpha$  ao conjunto  $P'$

        Fim de faça

    Fim de faça

Até não haver mais alterações em  $P'$

Fácil verificar que, ao fim deste algoritmo, todas as produções de  $P'$  já se encontram na FNG ■

Uma coisa interessante no algoritmo acima é que o mesmo não funciona para gramáticas com recursão esquerda (por que?).



# Exemplo

Seja a gramática LR(1) de expressões:

$$\begin{aligned}E &\rightarrow E + T \mid E - T \mid T \\T &\rightarrow T * F \mid T \div F \mid F \\F &\rightarrow (E) \mid \mathbf{n}\end{aligned}$$

$E$  é o símbolo inicial e  $\mathbf{n}$  é uma constante numérica qualquer. Com base na demonstração do teorema anterior, temos que

$$\begin{aligned}E &\rightarrow TR, & R &\rightarrow +TR \mid -TR \mid \varepsilon \\T &\rightarrow FQ, & Q &\rightarrow * FQ \mid \div FQ \mid \varepsilon \\F &\rightarrow (E) \mid \mathbf{n}\end{aligned}$$

Substituindo os prefixos com  $F$  nas demais produções, e observando que  $T \Rightarrow^* F$ , temos

$$\begin{aligned}E &\rightarrow (E)QR \mid \mathbf{n}QR, & R &\rightarrow +TR \mid -TR \mid \varepsilon \\T &\rightarrow (E)Q \mid \mathbf{n}Q, & Q &\rightarrow * FQ \mid \div FQ \mid \varepsilon \\F &\rightarrow (E) \mid \mathbf{n}\end{aligned}$$

Que já estão basicamente na forma normal de Greibach, a menos das produções épsilon. Contudo, se quisermos eliminar essas produções, iremos praticamente dobrar a complexidade espacial da gramática. O que importa é que todas as produções iniciam com terminais, o que era nosso objetivo neste exemplo.

# Gramática LL(1) – informal

A gramática de expressões do exercício anterior não precisava ficar explicitamente na forma normal de Greibach – seria suficiente eliminar as recursões esquerdas para atingirmos nosso próximo propósito. Tomemos a forma intermediária da gramática de expressões, onde não havia mais recursão esquerda:

$$\begin{aligned}E &\rightarrow TR, \\ R &\rightarrow +TR \mid -TR \mid \varepsilon \\ T &\rightarrow FQ, \\ Q &\rightarrow *FQ \mid \div FQ \mid \varepsilon \\ F &\rightarrow (E) \mid n\end{aligned}$$

Agora, observemos que a topologia de árvore de derivação, induzida pelas produções acima, é a mesma da hierarquia de chamada e implementação de funções C:  $E \rightarrow TR$  pode ser substituída por `void E(void){ T(); R(); }`. O mesmo se dá para as produções que partem de  $R$ :

```
1: void R(void) {
2:     if (lookahead == '+') {
3:         match ('+'); T(); R();
4:     } else if (lookahead == '-') {
5:         match ('-'); T(); R();
6:     } else {
7:         /*simula épsilon, caso venha lookahead == EOF ou lookahead == ')' */
8:         ;
9:     }
```

# O operador match

O operador match, que aparece na listagem anterior para simular as produções que partem de  $R$  tem o seguinte template:

```
1: void match(int expected_token) {  
2:     if (lookahead == expected_token) {  
3:         lookahead = gettoken(source);  
4:     } else {  
5:         fprintf(stderr, "token mismatch\n");  
6:         exit(ERRTOKEN);  
7:     }  
8: }
```

Então, conforme a listagem acima, match é a forma com que o parser solicita o próximo token, armazenador no operador lookahead. A parte correspondente ao erro não precisa ser do jeito que aparece nas linhas 5 e 6, e podem ser muito mais sofisticadas.

# A função FIRST

A gramática de expressões, conhecida na forma LL(1), foi escrita anteriormente como uma fase intermediária para passar a forma LR(1) para a forma normal de Greibach. Retomemos então a gramática LL(1) de expressões:

$$\begin{aligned} E &\rightarrow TR, & R &\rightarrow +TR \mid -TR \mid \varepsilon \\ T &\rightarrow FQ, & Q &\rightarrow^* FQ \mid \div FQ \mid \varepsilon \\ & & F &\rightarrow (E) \mid \mathbf{n} \end{aligned}$$

Confrontemos com a forma normal de Greibach:

$$\begin{aligned} E &\rightarrow (E)QR \mid \mathbf{n}QR, & R &\rightarrow +TR \mid -TR \mid \varepsilon \\ T &\rightarrow (E)Q \mid \mathbf{n}Q, & Q &\rightarrow^* FQ \mid \div FQ \mid \varepsilon \\ & & F &\rightarrow (E) \mid \mathbf{n} \end{aligned}$$

Observemos que, desta última forma, os primeiros tokens que podem ocorrer numa sentença válida gerada por esta gramática são '(' e **n**. Assim, uma sentença como "3 + 4" é válida, enquanto uma sentença como ")3 - + -4" não é válida, não somente pelos posteriores erros de sintaxe, pela repetição de operadores aditivos, mas por iniciar a mesma com ')'. Neste caso dizemos que '(' e **n** são membros de  $FIRST(E)$ . Adicionalmente, é fácil verificar que  $FIRST(E) = FIRST(T) = FIRST(F) = \{ (, \mathbf{n} \}$ . De modo análogo, as cadeias derivadas de  $R$  iniciam com os tokens em  $FIRST(R) = \{ \varepsilon, +, - \}$ .

Com base no que foi ensaiado acima,  $\varepsilon$  aparece como FIRST. Apesar de não fazer parte do alfabeto,  $\varepsilon$  é considerado um terminal, um terminal nulo. função FIRST:  $(\Sigma \cup N)^* \mapsto 2^{\Sigma \cup \{\varepsilon\}}$ , dada a gramática  $G = (N, \Sigma, S, P)$ , é a função que leva uma forma sentencial  $\alpha \in (\Sigma \cup N)^*$  numa coleção de tokens em  $2^{\Sigma \cup \{\varepsilon\}}$ , ou seja,  $FIRST(\alpha) \subseteq \Sigma \cup \{\varepsilon\}$ . No exemplo,  $\Sigma = \{ (, ), \mathbf{n}, +, -, *, \div \}$ .

# A função FIRST

Sumariamente, podemos definir a função FIRST conforme o seguinte algoritmo:

**Entrada:** uma gramática LL(1),  $G = (N, \Sigma, S, P)$ , um não símbolo de  $\Sigma \cup N \cup \{\varepsilon\}$  ou uma forma sentencial  $\alpha$  de  $(N \cup \Sigma)^*$ , onde  $A \rightarrow \alpha$  é uma produção de  $P$

**Saída:** Função  $\text{FIRST}(X)$

**Método:**

Inicialize  $\text{FIRST}(X) = \emptyset$

Se  $X \in N$  e  $X \rightarrow \varepsilon \in P$ , adicione  $\varepsilon$  a  $\text{FIRST}(X)$

Se  $X = \alpha_1 \alpha_2 \cdots \alpha_n$ , para algum  $A \rightarrow X \in P$ , ou  $X \rightarrow \alpha_1 \cdots \alpha_n$ , então

Se  $\alpha_i \Rightarrow^* \varepsilon$ , para  $i = 1, \dots, j - 1 < n$ , com  $\alpha_j \neq \varepsilon$ , inclua  $\text{FIRST}(\alpha_j)$  em  $\text{FIRST}(X)$

Senão

Se  $X \in \Sigma \cup \{\varepsilon\}$ , adicione  $X$  a  $\text{FIRST}(X)$

Fim de se

Retorne  $\text{FIRST}(X)$

Fim

# A função FOLLOW

Mais uma observação decorre do fato de que tanto  $R$  quanto  $Q$  retratam sentenças parciais, aparecendo como sufixo de  $T$  e como sufixo de  $F$  nas produções que partem de  $E$  e  $T$ , respectivamente.

Retomemos a versão LL(1) da gramática de expressões aritméticas:

$$\begin{aligned} E &\rightarrow TR, & R &\rightarrow +TR \mid -TR \mid \varepsilon \\ T &\rightarrow FQ, & Q &\rightarrow *FQ \mid \div FQ \mid \varepsilon \\ & & F &\rightarrow (E) \mid n \end{aligned}$$

Observemos que a primeira invocação do símbolo inicial  $E$  corresponde à raiz absoluta de uma árvore de derivação para, por exemplo, a expressão  $n * (n + n)$ . Neste caso, dizemos que  $E \Rightarrow^* n * (n + n)$ . Contudo, se quisermos enxergar a sentença produzida por uma gramática, incluindo o caráter de fim de sentença,  $\$$ , temos que, antes da invocação, o texto é simplesmente  $\$$ , o que significa que o texto é nulo, ou mais especificamente  $\$ = \varepsilon\$$ . Ao fim da derivação, o texto é visto, em baixo nível, como  $n * (n + n)\$$ . Fazendo uma redução na árvore, de baixo para cima, até chegarmos no  $E$ , isso equivale a ter escrito  $E\$$ . Ou seja, tudo que é derivado da invocação inicial do símbolo inicial é finalizado por  $\$$ . Então, dizemos que  $\$$  segue  $E$ , ou simplesmente que  $\$ \in \text{FOLLOW}(E)$ .

Raciocínio semelhante se dá com o que ocorre com  $(E)$ , quando  $E$  é invocado de dentro dos parênteses, vemos que tudo que for derivado de  $E$  é finalizado com  $)$ . Assim, como não outro símbolo seguindo  $E$  nas formas sentenciais da direita, temos finalmente  $\text{FOLLOW}(E) = \{\$, )\}$ .

Por outro lado, observe que  $E \Rightarrow TR\$$  ou  $E \Rightarrow TR)$  pode ocorrer ao longo da derivação de expressões. Então, tudo que for  $\text{FOLLOW}(E)$  servirá como  $\text{FOLLOW}(T)$  e é exatamente  $\text{FOLLOW}(R) = \text{FOLLOW}(E)$ , ou seja,  $\text{FOLLOW}(E) \subset \text{FOLLOW}(T)$ , o que ocorre quando  $R \Rightarrow \varepsilon$ . Contudo,  $T$  em  $TR$  vem seguido de tudo que estiver em  $\text{FIRST}(R)$  exceto  $\varepsilon$ . Então,  $+$  e  $-$  pertencem a  $\text{FOLLOW}(T)$ . Como a única aparição explícita de  $T$  ocorre somente em  $TR$ , concluímos que  $\text{FOLLOW}(T) = \text{FOLLOW}(R) \cup \text{FIRST}(R) - \{\varepsilon\} = \{+, -, ), \$\}$ .

# A função FOLLOW

Raciocínio semelhante ao que fizemos acima ocorre com  $Q$ . Como  $T \rightarrow FQ$ , tudo que segue  $T$  segue  $Q$ , e esta é a única aparição de  $Q$  partindo de outro não terminal, então  $\text{FOLLOW}(Q) = \text{FOLLOW}(T) = \{+, -, ), \$\}$ .

Generalizando, temos a definição de FOLLOW:  $N \mapsto 2^{\Sigma \cup \{\$ \}}$ , dada a gramática  $G = (N, \Sigma, S, P)$ , como a função que transforma um não terminal  $A \in N$  conforme o seguinte algoritmo:

1. Para cada não terminal  $A \in N$ , inicialize  $\text{FOLLOW}(A) = \emptyset$ ;
2. Se  $A = S$ , adicione  $\$$  a  $\text{FOLLOW}(A)$ ;
3. Se  $A$  aparece numa produção do tipo  $X \rightarrow \alpha A$ , inclua  $\text{FOLLOW}(X)$  a  $\text{FOLLOW}(A)$ ;
4. Se  $A$  aparece no lado direito de uma produção de  $P$ , na forma  $X \rightarrow \alpha A \beta$ , adicione  $\text{FIRST}(\beta) - \{\varepsilon\}$  a  $\text{FOLLOW}(A)$  – neste caso, se  $\varepsilon \in \text{FIRST}(\beta)$ , inclua  $\text{FOLLOW}(X)$  a  $\text{FOLLOW}(A)$ .

É imediato que  $\varepsilon \notin \text{FOLLOW}(X)$ .

# O que vem a ser afinal de contas uma gramática LL(1)?

Uma gramática  $G = (N, \Sigma, S, P)$ , livre do contexto, é dita LL(1) se é possível implementar um analisador sintático, codificado em C, que simule de forma única o processo de derivação dessa gramática, pela substituição de não terminais por procedimentos/funções C e seus terminais pela invocação do operador match desse terminal na forma de constante predefinida, sem retrocessos, ou nenhuma outra forma heurística, ou seja, de modo determinístico.

Em decorrência, uma gramática LL(1) tem que ter uma chave (token) única para indexar cada lado direito das regras de produção, ou seja,  $A \rightarrow \alpha \mid \beta$

Corresponde, virtualmente, ao trecho de código abstrato

```
1: void A(void) {
2:     if ( lookahead ∈ FIRST(α) ) {
3:         match(lookahead); /* continue denotando a forma sentencial α */
4:     else if ( lookahead ∈ FIRST(β) ) {
5:         match(lookahead); /* continue denotando a forma sentencial β */
6:     } else {
7:         /* codifique aqui a rotina de depuração de erro */
8:     }
9: }
```

Desta forma, com base no algoritmo acima, o lado direito  $\alpha$  é escolhido se e somente se o token previsto,  $lookahead \in FIRST(\alpha)$ , enquanto o lado direito  $\beta$  é escolhido se e somente se  $lookahead \in FIRST(\beta)$  o que só é possível se e somente se  $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$ . Isso torna determinística a escolha do lado direito das produções.

Por exemplo, a gramática  $S \rightarrow \mathbf{a} S \mid \mathbf{a} S \mathbf{b} S \mid \mathbf{c}$  não é LL(1). Por que?



# Ainda sobre LL(1)

Uma forma conveniente de enxergarmos o processo de análise sintática descendente preditiva recursiva como uma emulação do processo de derivação à esquerda é o seguinte.

Se  $S \rightarrow \alpha$  é uma produção inicial, então podemos idealizar a seguinte sucessão de descrições instantâneas:

$$.S \Rightarrow \alpha \vdash S \Rightarrow .\alpha \vdash^* S \Rightarrow \alpha. \equiv S. \Rightarrow \alpha$$

O ponto precedendo um símbolo significa a iminência de visitar o mesmo, sincronizado com *lookahead* – seja por chamada recursiva do procedimento/função associado, seja pela iminência de invocar o operador match do símbolo esperado em  $\text{FIRST}(\alpha)$ . O ponto após uma forma sentencial significa que o último componente, apêndice, foi visitado e isso significa o retorno de uma recursão ou saída do operador match.

Se tivermos uma produção com a seguinte configuração instantânea,

$$X \rightarrow \alpha.A$$

e  $\exists A \rightarrow \varepsilon \in P$ , então, a escolha  $X \Rightarrow \alpha.A \Rightarrow \alpha.\varepsilon \vdash \alpha\varepsilon. \equiv \alpha. \equiv X.$  é feita se e somente se *lookahead*, abstraindo o ponto prefixando  $A$ , não pertence a  $\text{FIRST}(A)$  mas pertence a  $\text{FOLLOW}(A)$ . O que não pode ocorrer é  $A \rightarrow \varepsilon$  para *lookahead*  $\in \text{FIRST}(A)$ , sabendo que *lookahead*  $\neq \varepsilon$ , o que tornaria possível ter *lookahead*  $\in \text{FIRST}(A) \cap \text{FOLLOW}(A) \neq \emptyset$  – em outras palavras, teríamos um *lookahead* sendo FIRST e FOLLOW ao mesmo tempo.

# Definição (finalmente) de gramática LL(1)

Uma gramática livre do contexto,  $G = (N, \Sigma, S, P)$ , é dita LL(1) se e somente se:

1. Todas as produções de  $P$  não possuem recursão esquerda – nem direta nem cíclica – em outras palavras,  $A \not\Rightarrow^* A\alpha$ ;
2. Uma produção do tipo  $A \rightarrow \alpha|\beta$  implica ter  $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$ ;
3. Se  $A \rightarrow \alpha|\varepsilon$ , com  $\alpha \neq \varepsilon$ , então a escolha  $A \rightarrow \varepsilon$  só possível quando *lookahead*  $\in \text{FOLLOW}(A)$ , com  $\text{FIRST}(A) \cap \text{FOLLOW}(A) = \emptyset$ .

## Exemplo de LL(1)

Surge o seguinte teorema:

**Teorema 2.** A gramática de expressões simplificadas,

$$G = (\{E, T, F, R, Q\}, \{\mathbf{n}, (, ), +, -, *, \div\}, E, P),$$

com

$$P = \{ E \rightarrow TR, R \rightarrow +TR, R \rightarrow -TR, R \rightarrow \varepsilon, T \rightarrow FQ, Q \rightarrow * FQ, Q \rightarrow \div FQ, Q \rightarrow \varepsilon, F \rightarrow (E), F \rightarrow \mathbf{n} \}$$

é uma gramática LL(1).

**Prova.** Segue imediatamente da definição de gramática LL(1) ■

# Exemplo de LL(1) – continuação

Tomemos a gramática LL(1) de expressões na sua forma mais abstrata,

$$\begin{aligned} E &\rightarrow TR \\ R &\rightarrow \oplus TR \mid \varepsilon \\ T &\rightarrow FQ \\ Q &\rightarrow \otimes FQ \mid \varepsilon \\ F &\rightarrow (E) \mid \mathbf{n} \end{aligned}$$

onde  $\oplus = ' + ' \mid ' - '$  e  $\otimes = ' * ' \mid \div$  abstraem os operadores aditivos e multiplicativos, respectivamente.

Neste caso, as possíveis configurações instantâneas para representar a análise sintática descendente preditiva recursiva são

$$.E \rightarrow TR \vdash E \rightarrow .TR \vdash^* E \rightarrow T.R \vdash^* E \rightarrow TR. \vdash E. \rightarrow TR$$

Ainda

$$\begin{aligned} .R \rightarrow \oplus TR \vdash R \rightarrow .\oplus TR \text{ (se } lookahead = \oplus) \vdash R \rightarrow \oplus.TR \vdash^* R \rightarrow \oplus T.R \vdash^* R \rightarrow \oplus TR. \vdash R. \rightarrow \oplus TR, \text{ ou} \\ R \rightarrow .\varepsilon \text{ (se } lookahead \notin FIRST(R) \wedge lookahead \in FOLLOW(R)) \vdash R \rightarrow \varepsilon. \vdash R. \rightarrow \varepsilon \end{aligned}$$

Mesmo raciocínio com  $T \rightarrow FQ$  e  $Q \rightarrow \otimes FQ \mid \varepsilon$ , mas dando ênfase às produções finais

$$\begin{aligned} .F \rightarrow (E) \vdash F \rightarrow .(E) \vdash F \rightarrow (.E) \vdash^* F \rightarrow (E.) \vdash F \rightarrow (E). \vdash F. \rightarrow (E) \\ .F \rightarrow \mathbf{n} \vdash F \rightarrow .\mathbf{n} \vdash F \rightarrow \mathbf{n}. \vdash F. \rightarrow \mathbf{n} \end{aligned}$$

# Exemplo de LL(1) – continuação

Onde as funções FIRST e FOLLOW foram computadas como

$$\begin{aligned}\text{FIRST}(E) &= \text{FIRST}(T) = \text{FIRST}(F) = \{n, ()\} \\ \text{FIRST}(R) &= \{\varepsilon, \oplus\} \\ \text{FIRST}(Q) &= \{\varepsilon, \otimes\}\end{aligned}$$

$\text{FOLLOW}(E) = \{\$, )\}$ , visto que podem surgir as situações  $E.\$$  e  $E.)$

$\text{FOLLOW}(T) = \text{FOLLOW}(E) \cup \text{FIRST}(R) - \{\varepsilon\} = \{\$, ), \oplus\}$ , visto que em na configuração instantânea  $T.R$ , a produção  $R \rightarrow \varepsilon$  ocorre se *lookahead*  $\in \text{FOLLOW}(R)$ , ou então  $T$  é seguido de  $\oplus$ , caso *lookahead*  $= \oplus$

$\text{FOLLOW}(R) = \text{FOLLOW}(E) = \{\$, )\}$ , uma vez que  $E \rightarrow TR$ . é seguido de tudo que segue  $E$

$\text{FOLLOW}(Q) = \text{FOLLOW}(T) = \{\$, ), \oplus\}$ , uma vez que  $T \rightarrow FQ$ . é seguido de tudo que segue  $T$

Por último,  $\text{FOLLOW}(F) = \text{FOLLOW}(T) \cup \text{FIRST}(Q) - \{\varepsilon\}$ , visto que em na configuração instantânea  $F.Q$ , a produção  $Q \rightarrow \varepsilon$  ocorre se *lookahead*  $\in \text{FOLLOW}(Q)$ , ou então  $F$  é seguido de  $\otimes$ , caso *lookahead*  $= \otimes$

**É sempre conveniente justificar cada passo de obtenção das funções recursivas acima, principalmente em caso de provas ou de outras atividades para nota.**

# Exemplo de LL(1) – continuação

Continuando o exemplo, suponhamos agora que queremos derivar  $3 + 4$ , que será visto pelo analisador sintático como  $n \oplus n$ , utilizando a notação ponto:

$$\begin{aligned} & .E \Rightarrow TR \vdash E \Rightarrow .TR \Rightarrow .FQR \text{ (lookahead} = n) \\ & \Rightarrow .nQR \vdash n.QR \text{ (lookahead} = \oplus) \Rightarrow n.\varepsilon R \vdash n.R \\ & \Rightarrow n.\oplus TR \vdash n \oplus .TR \Rightarrow n \oplus .FQR \text{ (lookahead} = n) \\ & \Rightarrow n \oplus .nQR \vdash n \oplus n.QR \text{ (lookahead} = \$) \\ & \Rightarrow n \oplus n.\varepsilon R \vdash n \oplus n.R \Rightarrow n \oplus n.\varepsilon \vdash n \oplus n. \\ & \vdash E. \Rightarrow TR \end{aligned}$$

Este último passo é só para dizer que retornou da recursão inicial, da derivação  $.E \Rightarrow TR$

# Como converter uma gramática livre do contexto para LL(1)?

Pelo que pudemos perceber da explanação anterior, uma gramática é LL(1) se não tem recursão esquerda e o lado direito de cada produção, digamos, a forma sentencial  $\alpha$  em  $A \rightarrow \alpha$ , é única em relação ao token lido antecipadamente em *lookahead*, combinado com o símbolo não terminal,  $A$ , da esquerda. Assim, precisamos encontrar uma gramática cujo lado direito das produções  $A \rightarrow \alpha$  seja uma função  $\alpha = f(A, \tau)$ , onde  $\tau$  é o token carregado em *lookahead*. Assim, podemos fazer as seguintes prescrições:

1. Elimine as recursões esquerdas das produções, do mesmo modo que fizemos no início desta aula;
2. Faça a fatoração esquerda nas formas sentenciais direitas das produções – por exemplo, se tivermos  $S \rightarrow aS \mid aSbS \mid c$ , podemos rescrever o lado direito como se fosse uma expressão regular, obtendo  $aS \mid aSbS \mid c = aS(\varepsilon \mid bS) \mid c$ . Como assumimos como dogma que não iremos misturar expressão regular com gramática, redefinimos o termo entre parênteses como  $R \rightarrow bS \mid \varepsilon$ . Deste modo, a gramática resultante é  $S \rightarrow aSR \mid c, R \rightarrow bS \mid \varepsilon$ ;
3. Estude FIRST e FOLLOW para saber o conjunto de tokens necessário para selecionar os lados direitos e para saber quando algum não terminal deve se anular – por exemplo,  $R \rightarrow \varepsilon$ , do exemplo dado no Item 2, é necessário sempre que o token  $\tau$  armazenado em *lookahead* estiver em  $\text{FOLLOW}(R) = \{\$ \}$

Observe que a condição do Item 3 remove a ambiguidade que existe na gramática, uma vez que  $R$  nunca é zerado quando *lookahead* =  $b$ . Por exemplo, derivando  $S \Rightarrow aSR \vdash a.SR \Rightarrow a.aSRR \vdash aa.SRR \vdash^* aaS.RR$  – neste momento, *lookahead* pode ser  $\$$  ou pode ser  $b$ . No primeiro caso, temos que  $.RR \Rightarrow \varepsilon.R \Rightarrow \varepsilon\varepsilon$ . No segundo caso, quando *lookahead* =  $b$ , só podemos ter o  $R$  prefixado pelo ponto, ficando  $.RR \Rightarrow b.R$ , o que impede que o primeiro  $R$  seja anulado na presença de *lookahead* =  $b$ .

# Mais um exemplo de analisador descendente recursivo

A gramática

$$\begin{aligned} S &\rightarrow aSR \mid c, \\ R &\rightarrow bS \mid \varepsilon \end{aligned}$$

pode ser facilmente transcrita para um analisador sintático recursivo preditivo:

```
1: void S(void) {
2:     if (lookahead == a) {
3:         match(a); S(); R();
4:     } else {
5:         match(c);
6:     }
7: }
```

```
1: void R(void) {
2:     if (lookahead == b) {
3:         match(b); S();
4:     } else {
5:         match(EOF); /* gettoken embutido em
6:         match vai retornar EOF se chamado
7:         novamente */
8:     }
9: }
```



# Testando o exemplo anterior

Para entender o funcionamento de ambas rotinas simulando a referida gramática, tomemos a sequência de tokens  $|a|a|c|$ , onde  $\$$  simboliza a constante EOF. Para simplificar, supomos que os espaços foram removidos. A sequência de chamadas é simulada por

$$.S \Rightarrow .aSR \vdash a.SR \Rightarrow a.aSRR \vdash aa.SRR \Rightarrow aa.cRR \vdash aac.RR$$

Neste momento, onde o ponto está prefixando o  $R$  mais à esquerda da dupla  $RR$ , *lookahead* já é EOF. Isso corresponde à linha 5 do procedimento  $R()$  da listagem anterior.

Então  $\text{match}(\text{EOF})$  vai ser bem sucedido e vai invocar  $\text{gettoken}$ , que continuará retornando EOF.

Acontece que a linha 5 reproduz  $R \rightarrow \varepsilon$ . Logo, teremos

$$aac.RR \Rightarrow aac.\varepsilon R \vdash aac\varepsilon.R \Rightarrow aac\varepsilon.\varepsilon \vdash aac\varepsilon\varepsilon. \equiv aac.$$

que significa que acabou a análise com sucesso.

# Questões valendo 2 horas – vamos ver se prestaram atenção à aula assíncrona?

1. Explique com poucas palavras, e com suas palavras, o que vem a ser a forma normal de Greibach.
2. Por que o algoritmo, utilizado para concluir a demonstração do Teorema de Greibach nem sempre funciona com recursão esquerda? *Dica: teste com a gramática LR(1) de expressões.*
3. Qual a vantagem de se ter a gramática de expressões na forma LL(1) do ponto de vista de implementação simplificada de analisadores sintáticos?
4. Porque  $S \rightarrow \mathbf{a}S \mid \mathbf{a}S\mathbf{b}S \mid \varepsilon$  não é LL(1)?
5. A gramática  $S \rightarrow \mathbf{a}SR, R \rightarrow \mathbf{b}S \mid \varepsilon$ , com  $S$  sendo o símbolo inicial, é LL(1)? Justifique.