

## Лекція №3 Граматики.

### Формальне визначення граматики. Типи граматик і їх властивості.

Для нас найбільший інтерес представляє одна з систем генерації мов - граматики. Поняття граматики спочатку було формалізовано лінгвістами при вивченні природних мов. Передбачалося, що це може допомогти при їх автоматичній трансляції. Однак, найкращі результати в цьому напрямку досягнуто при описі не природних мов, а мов програмування. Прикладом може служити спосіб опису синтаксису мов програмування за допомогою БНФ - форми Бекуса-Наура.

Визначення. Граматика - це четвірка  $G = (N, T, P, S)$ , де

$N$  - алфавіт нетермінальних символів;

$T$  - алфавіт термінальних символів,  $NT = \emptyset$ ;

$P$  - кінцева множина правил виду, де  $(NT)^* N (NT)^*, (NT)^*$ ;

$S$  - початковий символ (або аксіома) граматики.

Ми будемо використовувати великі латинські літери для позначення нетермінальних символів, малі латинські букви з початку алфавіту для позначення термінальних символів, малі латинські букви з кінця алфавіту для позначення ланцюжків з  $T^*$ , нарешті, малі грецькі літери для позначення ланцюжків з  $(NT)^*$ .

Будемо використовувати також скорочений запис  $A_1 | A_2 | \dots | A_n$  для позначення групи правил  $A_1, A_2, \dots, A_n$ .

Визначимо на множині  $(NT)^*$  бінарне відношення виводимості таким чином: якщо  $P$ , то для всіх,  $(NT)^*$ . Якщо  $1 \Rightarrow 2$ , то говорять, що ланцюжок  $2$  безпосередньо виведена з  $1$ .

Ми будемо використовувати також рефлексивно-транзитивне і транзитивне замикання відношення, а також його ступінь  $k \geq 0$  (позначаються відповідно  $*$ ,  $+$  і  $k$ ). Якщо  $1 \Rightarrow^k 2$  ( $k \geq 0$ ), то говорять, що ланцюжок  $2$  виведена (нетривіально виведена, виведена за  $k$  кроків) з  $1$ .

Якщо  $k \geq 0$ , то існує послідовність кроків

де  $i = 0, 1, 2, \dots, k$ . Послідовність ланцюжків  $0, 1, 2, \dots, k$  в цьому випадку називають виведенням з.

Сентенціальний формою граматики  $G$  називається ланцюжок, що виводиться з її початкового символу.

Мовою, породжуваною граматиною  $G$  (позначається  $L(G)$ ), називається множина всіх її термінальних сентенціальних форм, тобто

Граматики  $G_1$  і  $G_2$  називаються еквівалентними, якщо вони породжують одну й ту саму мову, тобто  $L(G_1) = L(G_2)$ .

Приклад 2.5. Граматика  $G = (\{S, B, C\}, \{a, b, c\}, P, S)$ , де

$P = \{S \rightarrow aSBC, S \rightarrow aBC, CB \rightarrow BC, aB \rightarrow ab, bB \rightarrow bb, bC \rightarrow bc, cC \rightarrow cc\}$ , породжує мову  $L(G) = \{anbncn \mid n > 0\}$ .

Дійсно, застосовуємо  $n-1$  раз правило 1 і отримуємо  $a^{n-1}S(BC)^{n-1}$ , потім один раз правило 2 і отримуємо  $a^n(BC)^n$ , потім  $n(n-1)/2$  раз правило 3 і отримуємо  $a^nB^nC^n$ .

Потім використовуємо правило 4 і отримуємо  $anbB^{n-1}C^n$ . Потім застосовуємо  $n-1$  раз правило 5 і отримуємо  $anb^nC^n$ . Потім застосовуємо правило 6 і  $n-1$  раз правило 7 і отримуємо  $anbnc^n$ . Можна показати, що мова  $L(G)$  складається з ланцюжків тільки такого виду.

Приклад 2.6. Розглянемо граматику  $G = (\{S\}, \{0, 1\}, \{S \rightarrow 0S1, S \rightarrow 01\}, S)$ . Легко бачити, що ланцюжок  $000111 \in L(G)$ , так як існує висновок

Неважко показати, що граматика породжує мову  $L(G) = \{0^n1^n \mid n > 0\}$ .

Приклад 2.7. Розглянемо граматику  $G = (\{S, A\}, \{0, 1\}, \{S \rightarrow 0S, S \rightarrow 0A, A \rightarrow 1A, A \rightarrow 1\}, S)$ . Неважко показати, що граматика породжує мову  $L(G) = \{0^n1^m \mid n, m > 0\}$ .

## Типи граматик і їх властивості

Розглянемо класифікацію граматик (запропоновану Н. Хомського), засновану на вигляді їх правил.

Визначення. Нехай дана граматика  $G = (N, T, P, S)$ . Тоді

якщо правила граматики не задовольняють жодним обмеженням, то її називають граматикою типу 0, або граматикою без обмежень.

якщо

кожне правило граматики, крім  $S \rightarrow \epsilon$ , має вигляд  $de || || |$ , і

в тому випадку, коли  $S \in P$ , символ  $S$  не зустрічається в правих частинах правил,

то граматика називається граматикою типу 1, або нескороченою.

якщо кожне правило граматики має вигляд  $A \rightarrow \alpha$ , де  $A \in N$ ,  $(NT)^*$ , то її називають граматикою типу 2, або контекстно-вільною (КС-граматикою).

якщо кожне правило граматики має вигляд або  $A \rightarrow xB$ , або  $A \rightarrow x$ , де  $A, B \in N$ ,  $x \in T^*$  то її називають граматикою типу 3, або праволінійною.

Легко бачити, що граматика в прикладі 2.5 - нескороченою, в прикладі 2.6 - контекстно-вільна, в прикладі 2.7 - праволінійна.

Мова, породжується граматикою типу  $i$ , називається мовою типу  $i$ . Мова типу 0 називається також мовою без обмежень, мова типу 1 - контекстно-залежним (КЗ), мова типу 2 - контекстно-вільним (КС), мова типу 3 - праволінійним.

Теорема 2.1. Кожна контекстно-вільна мова може бути породжена неукорачивающей контекстно-вільною граматикою.

Доказ. Нехай  $L$  - контекстно-вільна мова. Тоді існує контекстно-вільна граматика  $G = (N, T, P, S)$ , що породжує  $L$ .

Побудуємо нову граматика  $G' = (N', T, P', S')$  таким чином:

1. Якщо в  $P$  є правило виду  $A \rightarrow B_1 B_2 \dots B_k$ , де  $k \geq 0$ ,  $B_i \neq \epsilon$  для  $1 \leq i \leq k$ , і ні з одного ланцюжка  $B_j$  ( $0 \leq j < k$ ) не виводиться  $\epsilon$ , то включити в  $P'$  всі правила (крім  $A \rightarrow \epsilon$ ) виду

де  $X_i$  - це або  $B_i$ , або  $\epsilon$ .

2. Якщо  $S \neq \epsilon$ , то включити в  $P'$  правила  $S' \rightarrow S$ ,  $S' \rightarrow \epsilon$  і покласти  $N' = N \cup \{S'\}$ . В іншому випадку покласти  $N' = N$  і  $S' = S$ .

Чи породжує граматика порожню ланцюжок можна встановити наступним простим алгоритмом:

Крок 1. Будуємо множину  $N_0 = N \mid N \rightarrow e$

Крок 2. Будуємо множину  $N_i = N \mid N \rightarrow \alpha \in N_{i-1}^*$

Крок 3. Якщо  $N_i = N_{i-1}$ , перейти до кроку 4, інакше крок 2.

Крок 4. Якщо  $S \rightarrow N_i$  значить  $S \rightarrow e$ /

Легко бачити, що  $G'$  - неукорачивающая граматика. Можна показати по індукції, що  $L(G') = L(G)$ . \_\_

Нехай  $K_i$  - клас всіх мов типу  $i$ . Доведено, що справедливо наступне (строге) включення:

$$K_3 \subset K_2 \subset K_1 \subset K_0.$$

Зауважимо, що якщо мова породжується деякої граматикою, це не означає, що він не може бути породжений граматикою з більш сильними обмеженнями на правила. Наведений нижче приклад ілюструє цей факт.

Приклад 2.8. Розглянемо граматика  $G = (\{S, A, B\}, \{0, 1\}, \{S \rightarrow AB, A \rightarrow 0A, A \rightarrow 0, B \rightarrow 1B, B \rightarrow 1\}, S)$ . Ця граматика є контекстно-вільною. Легко показати, що  $L(G) = \{0^n 1^m \mid n, m > 0\}$ .

Однак, в прикладі 2.7 наведено праволінійная граматика, що породжує ту ж мову.

Покажемо що існує алгоритм, що дозволяє для довільного КЗ-мови  $L$  в алфавіті  $T$ , і довільної ланцюжка  $w \in T^*$  визначити, чи належить  $w$  мові  $L$ .

Теорема 2.2. Кожен контекстно-залежний мову є рекурсивним мовою.

Доказ. Нехай  $L$  - контекстно-залежний мову. Тоді існує деяка неукорачивающая граматика  $G = (N, T, P, S)$ , що породжує  $L$ .

Нехай  $w \in T^*$  і  $|w| = n$ . Якщо  $n = 0$ , тобто  $w = e$ , то приналежність  $w \in L$  перевіряється тривіальним чином. Так що будемо припускати, що  $n > 0$ .

Визначимо множина  $T_m$  як множина рядків  $u \in (N \cup T)^+$  довжини не більше  $n$  таких, що висновок  $S \Rightarrow^* u$  має не більше  $m$  кроків. Ясно, що  $T_0 = \{S\}$ .

Легко показати, що  $T_m$  можна отримати з  $T_{m-1}$  переглядаючи, які рядки з довжиною, меншою або рівною  $n$  можна вивести з рядків з  $T_{m-1}$  застосуванням одного правила, тобто

$$T_m = T_{m-1} \cup \{u \mid v \Rightarrow u \text{ для некоторого } v \in T_{m-1}, \text{ где } |u| \leq n\}.$$

Якщо  $S \Rightarrow^* u$  и  $|u| \leq n$ , то  $u \in T_m$  для деякого  $m$ . Якщо з  $S$  не виводиться  $u$  або  $|u| > n$ , то  $u$  не належить  $T_m$  ні для якого  $m$ .

Очевидно, що  $T_m \supseteq T_{m-1}$  для всіх  $m \geq 1$ . Оскільки  $T_m$  залежить тільки від  $m \geq 1$ , якщо  $T_m = T_{m-1}$ , то  $T_m = T_{m+1} = T_{m+2} = \dots$ . Процедура буде обчислювати  $T_1, T_2, T_3, \dots$  поки для деякого  $m$  не виявиться  $T_m = T_{m-1}$ . Якщо  $w$  не належить  $T_m$ , то не належить і  $L(G)$ , оскільки для  $j > m$  виконано  $T_j = T_m$ . Якщо  $w \in T_m$ , то  $S \Rightarrow^* w$ .

Покажемо, що існує таке  $m$ , що  $T_m = T_{m-1}$ . Оскільки для кожного  $i \geq 1$  справедливо  $T_i \supseteq T_{i-1}$ , то якщо  $T_i \neq T_{i-1}$ , то число елементів в  $T_i$  принаймні на 1 більше, ніж в  $T_{i-1}$ . Нехай  $|N \cup T| = k$ . Тоді число рядків в  $(N \cup T)^+$  довжини меншою або рівною  $n$  одно  $k + k^2 + \dots + k^n \leq nk^n$ . Тільки ці рядки можуть бути в будь-якому  $T_i$ . Значить,  $T_m = T_{m-1}$  для деякого  $m \leq nk^n$ . Таким чином, процедура, що обчислює  $T_i$  для всіх  $i \geq 1$  до тих пір, поки не будуть знайдені два рівних множини, гарантовано закінчується, значить, це алгоритм.