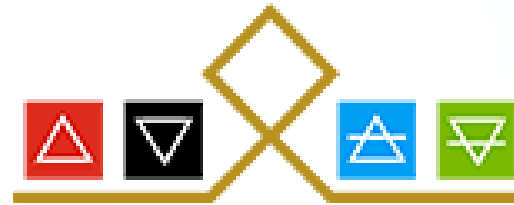




НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ  
БІОРЕСУРСІВ І ПРИРОДОКОРИСТУВАННЯ  
УКРАЇНИ



# Теорія розпізнавання образів та класифікації в системах штучного інтелекту

*Тема №7. Семантичний аналіз текстів формальної  
мови*

Київ - 2025

# Зміст

- 1. Семантичний аналіз.*
- 2. Семантичний аналізатор.*

## Список використаних джерел

1. Серебряков - Мови програмування:  
<http://infonet.cherepovets.ru/citforum/programming/theory/serebryakov>
2. Вільна енциклопедія - Вікіпедія <http://ru.wikipedia.org/wiki/%D0%A2%D1%80%D0%B0%D0%BD%D1%81%D0%BB%D1%8F%D1%82%D0%BE%D1%80>

На етапі контекстного аналізу виявляються залежності між частинами програми, які не можуть бути описані синтаксисом контекстно-вільної мови. Полягає в основному зв'язку, а саме - «опис-використання», а саме: аналіз типів об'єктів, аналіз областей видимості, відповідність параметрів, мітки та інші. У процесі контекстного аналізу таблиці об'єктів поповнюються інформацією про описи (властивості) об'єктів.

Основним формалізмом, що використовується при контекстному аналізі, є апарат атрибутних граматик.

Результатом контекстного аналізу є атрибутовані дерево програми. Інформація<sup>I</sup> про об'єкти може бути як розосереджена в самому дереві, так і зосереджена в окремих тубличних об'єктах.

У процесі контекстного аналізу також можуть бути виявлені помилки, пов'язані з неправильним використанням об'єктів.

Потім програма може бути переведена у внутрішнє представлення. Це робиться для цілей оптимізації та/або зручності генерації коду.

Ще однією метою перетворення програми у внутрішнє представлення є використання переносимого компілятора. В такому разі, тільки остання фаза (генерація коду) є машинно-залежною. В якості внутрішнього подання може використовуватися префіксний або постфіксний запис, орієнтований граф, трійки, четвірки та інші.

## *Основні принципи роботи синтаксичного аналізатора*

Синтаксичний аналізатор (синтаксичний розбір) - це частина компілятора, яка відповідає за виявлення основних синтаксичних конструкцій вхідної мови. У завдання синтаксичного аналізу входить: знайти і виділити основні синтаксичні конструкції в тексті вхідної програми, встановити тип і перевірити правильність кожної синтаксичної конструкції і, нарешті, представити синтаксичні конструкції у вигляді, зручному для подальшої генерації тексту результуючої програми.

В основі синтаксичного аналізатора лежить розпознаватель тексту вхідної програми на основі граматики вхідної мови. Як правило, синтаксичні конструкції мов програмування можуть бути описані за допомогою КВ-грамматик, рідше зустрічаються мови, які, можуть бути описані за допомогою регулярних грамматик. Найчастіше регулярні граматики застосовні до мов асемблера, а мови високого рівня побудовані на основі синтаксису КВ-мов. Розпізнавач дає відповідь на питання про те, належить чи не належить ланцюжок вхідних символів заданому в мові.

Як і у випадку лексичного аналізу, завдання синтаксичного розбору не обмежилося лише перевіркою приналежності ланцюжка заданому мови. Необхідно виконати всі перераховані вище завдання, які повинен вирішити синтаксичний аналізатор. У такому варіанті аналізатор вже не є різновидом МП-автомата – його функції можна трактувати ширше. Синтаксичний аналізатор повинен мати якиюсь вихідну мову, за допомогою якої він передає наступним фазам компіляції не тільки інформацію про знайдені і розібрані синтаксичні структури.

У такому випадку він є перетворювачем з магазинною пам'яттю - МП-перетворювачем.

Синтаксичний розбір - це основна частина компілятора на етапі аналізу.

Усі завдання з перевірки синтаксису вхідного мови можуть бути вирішені на етапі синтаксичного розбору.

Сканер тільки дозволяє позбавити складний за структурою синтаксичний аналізатор від рішення примітивних завдань з виявлення та запам'ятовування лексем вихідної програми.

Виходом лексичного аналізатора є таблиця лексем (або ланцюжок лексем). Ця таблиця утворює вхід синтаксичного аналізатора, який досліджує тільки один компонент кожної лексеми - її тип. Решта інформації про лексемах використовується на більш наступних етапах компіляції при семантичному аналізі, підготовці до генерації і генерації коду результуючої програми.

Синтаксичний аналіз (розбір) - це [процес](#), в якому досліджується таблиця лексем і встановлюється, чи задовольняє вона структурним умов, явно сформульованим і визначених синтаксисом мови.

Синтаксичний аналізатор сприймає вихід лексичного аналізатора і розбирає його відповідно до граматики вхідної мови.

## *Дерево розбору. Перетворення дерева розбору в дерево операцій*

Результатом роботи розпізнавача КС-граматики вхідної мови є послідовність правил граматики, застосованих для побудови вхідний ланцюжка. За знайденою послідовністю, знаючи тип розпізнавача, можна побудувати ланцюжок виводу або дерево виводу. У цьому випадку дерево виведення виступає в якості дерева синтаксичного розбору і являє собою результат роботи синтаксичного аналізатора в компіляторі.

Ні ланцюжок виводу, ні дерево синтаксичного розбору не є метою роботи компілятора. Дерево висновку містить масу надлишкової інформації, яка для подальшої роботи компілятора не потрібна. Ця інформація включає в себе всі нетермінальні сисволи, які містяться і вузлах дерева, Після побудови дерева нетермінальні символи не несуть ніякого смислового навантаження і не використовуються.

Для повного уявлення про тип і структурі знайденої і розібраної синтаксичної конструкції вхідного мови достатньо знати послідовність номерів правил граматики, що використовуються для її побудови. Форма подання цієї достатньої інформації може бути різною як залежно від реалізації самого компілятора, так від фази компіляції. Ця форма напивається внутрішнім поданням програми.

У синтаксичному дереві внутрішні вузли (вершини) відповідають операціям, а листя представляють собою операнди. Як правило, листя синтаксичного дерева співпадають з записами в таблиці ідентифікаторів.

Структура синтаксичного дерева відображає синтаксису мови програмування, на якому нашкапа похідна програма.

Синтаксичні дерева можуть бути побудовані компілятором для будь-якої частини вхідної програми. Не завжди синтаксичному дереву повинен відповідати фрагмент коду результуючої програми, наприклад, можлива побудова синтаксичних дерев для декларативної частини мови. У цьому випадку операції, наявні в дереві, не вимагають породження об'єктного коду, але несуть інформацію про дії, які повинен виконати сам компілятор над відповідними елементами. У випадку, коли синтаксичному дереву відповідає деяка послідовність операцій, що тягне породження фрагмента об'єктного коду, говорять про дерево операцій.

Дерево операцій можна безпосередньо побудувати з дерева виведення, породженого синтаксичним аналізатором. Для цього достатньо виключити з дерева виведення ланцюжка нетермінальних символів, а також вузли, що не несуть семантичного навантаження при генерації коду.

Прикладом таких вузлів можуть бути різні дужки, які змінюють порядок виконання операції і операторів, але після побудови дерева ніякого смислового навантаження не несуть, гак яким не відповідає ніякий об'єктним код.

## *Алгоритм перетворення дерева семантичного розбору і дерево операції*

*Крок 1.* Якщо у дереві більше не міститься вузлів, помічених нетермінальним символами, то виконання алгоритму завершено; інакше - перейти до кроку 2.

*Крок. 2.* Вибрати крайній лівий вузол лерена, що позначений нетермінальним символом граматика і зробити його поточним. Перейти до кроку 3.

*Крок 3.* Якщо поточний вузол має тільки один нижє розташований вузол, то поточний вузол необхідно видалити з дерена, а пов'язаний з ним вузол приєднати до вузла вищого рівня (виключити з лерена ланцюжок) і повернутися до кроку 1; інакше - перейти до кроку 4.

*Крок 4.* Якщо поточний вузол має нижєлежачий вузол (лист дерева), позначений термінальним символом, який не несе семантичного навантаження, тоді цей лист потрібно видалити з дерева і повернутися до кроку 3, інакше - перейти до кроку 5.

*Крок 5.* Якщо поточний вузол має один нижєлежачий вузол (лист дерева), позначений термінальним символом, що позначає знак операції, а інші вузли позначені як операнди, то лист, позначений знаком операції, треба видалити з дерева, поточний вузол позначити цим знаком операції і перейти до кроку 1; інакше - перейти до кроку 6.

*Крок 6.* Якщо серед нижчих вузлів для поточного вузла є вузли, помічені нетермінальним символами граматика, то необхідно вибрати крайній лівий серед цих вузлів, зробити його поточним і перейти до кроку 3: інакше - виконання алгоритму завершено.

# 1. Семантичний аналіз

КВ-граматики, за допомогою яких описують синтаксис мов програмування, не дозволяють задавати контекстні умови (КУ), що є в будь-якій мові. Перевірку КУ називають **семантичним аналізом**.

Найчастіше зустрічаються контекстні умови:

- кожен ідентифікатор, що використовується в програмі, повинен бути описаний, але не більше одного разу на одній зоні опису;
- при виклику функції число та тип фактичних параметрів повинні відповідати числу та типам формальних параметрів;
- зазвичай у мові накладаються обмеження на:
  - типи операндів будь-якої операції, визначеної у цій мові;
  - типи лівої та правої частин в операторі присвоєння;
  - Тип параметра циклу;
  - тип умови в операторах циклу та умовному операторі тощо.

Перевірка КУ проводиться в залежності від того, як синтаксичний аналізатор розпізнає конструкцію, на компоненти якої накладені деякі обмеження, тобто. на етапі синтаксичного аналізу повинні виконуватись деякі додаткові дії, які здійснюють семантичний контроль.

## Вставка дій у КС-граматику

Якщо для синтаксичного аналізу використовується **метод рекурсивного спуску (РС)**, то контролю КУ в тіла процедур РС-метода необхідно вставити виклики додаткових " семантичних " процедур (семантичні дії).

Спочатку семантичні дії вставляються в синтаксичні правила, а потім за цими розширеними правилами будуються процедури РС-методу.

Щоб відрізнити виклики семантичних процедур від інших символів граматики, вони полягають у кутових дужках.

Фактично розширюється поняття КВ-грамматики.

Нехай у граматиці є правило

$$A \rightarrow a \langle D1 \rangle B \langle D1; D2 \rangle \mid b \langle C \langle D3 \rangle \rangle,$$
 де

$A, B, C \in N; a, b \in T; \langle Di \rangle$  - є виклик процедури  $Di, i = 1, 2, 3$ .

За таким правилом граматики процедуру для РС-методу буде такою:

```
void A ( ) {
if (c == 'a') { gc(); D1(); B(); D1(); D2(); }
else
if (c == 'b') { gc(); C(); D3(); }
else throw c;
}
```

## Приклад

Написати граматику, яка породжує мову

$L = \{ \alpha \in (0,1)^+ \mid \alpha \text{ містить рівну кількість } 0 \text{ і } 1 \}$ .

Вирішити це завдання можна

- чисто синтаксичними засобами - описати ланцюжки, які мають необхідну властивість;

- за допомогою синтаксичних правил описати довільні ланцюжки з 0 та 1, а потім вставити дії для відбору ланцюжків з рівною кількістю 0 та 1.

$S \rightarrow \langle k_0 = k_1 = 0; \rangle A \perp$

$A \rightarrow 0 \langle k_0++; \rangle B \mid 1 \langle k_1++; \rangle B$

$B \rightarrow A \mid \epsilon \langle \text{if } (k_0 \neq k_1) \text{ throw "ERROR !!!"; } \rangle$

Контекстні умови, виконання яких треба контролювати у програмах на М-мові, такі:

- Будь-яке ім'я, що використовується в програмі, має бути описано і лише один раз.
- В операторі присвоювання типи змінної та виразу повинні збігатися.
- В умовному операторі та в операторі циклу як умови можливе лише логічне вираження.
- Операнди операції відносини мають бути цілими.
- Тип вираження та сумісність типів операндів у виразі визначаються за звичайними правилами (як у Паскалі).

Для перевірки КУ М-мови в синтаксичні правила граматики вставимо виклики процедур, які здійснюють необхідний контроль, а потім перенесемо в процедури рекурсивного спуску.

У синтаксичні правила для описів потрібно вставити дії, за допомогою яких можна запам'ятати тип змінних та контролювати єдиність їх опису.

$i$ -ий рядок таблиці TID відповідає ідентифікатору-лексемі виду (LEX\_ID, $i$ ).

Лексичний аналізатор заповнив поле **name**; значення полів **declare** і **type** заповнюватимемо на етапі семантичного аналізу.

Розділ описів має вигляд

**D** → **I** { ,**I** } : [ **int** | **bool** ],

тобто. імені типу (**int** або **bool**) передуює список ідентифікаторів.

Ці ідентифікатори (номери відповідних рядків таблиці TID) треба запам'ятати, наприклад, у стеку цілих чисел

*Stack* < *int*, 100 > *st\_int* ,

а коли буде проаналізовано ім'я типу, треба заповнити поля **declare** і **type** у відповідних рядках.

Функція *void Parser::dec(type\_of\_lextype)*:

- ✓ зчитує зі стека номера рядків таблиці TID,
- ✓ заносить у них інформацію про тип відповідних змінних та про наявність їх описів та
- ✓ контролює повторне опис змінних.

```

parser::dec ( type_of_lex type ) {
    int i;
    while (! st_int.is_empty ( )) {
        i = st_int.pop ();
        if ( TID [ i ].get_declare () throw "twice";
        else {
            TID [ i ].put_declare ();
            TID [ i ].put_type (type);
        }
    }
}

```

З урахуванням наявних функцій правило виведення з діями обробки описів буде ТАКИМ:

```

D → < st_int.reset () > I < st_int.push (c_val) >
{ , I < st_int.push (c_val) > } :
[ int < dec (LEX_INT) > | bool < dec (LEX_BOOL) > ]

```

## Контроль контекстних умов у виразі

Типи операндів і позначення операцій зберігатимемо в стеку *Stack*`<type_of_lex, 100>st_lex`.

Якщо у виразі зустрічається лексема-ціле число чи логічні константи *true* або *false* то відповідний тип відразу заноситься в стек.

Якщо операнд - лексема-змінна, необхідно перевірити, чи описана вона; якщо описана, її тип треба занести в стек.

Ці дії виконуються за допомогою функції `check_id`:

```
void parser::check_id ( ) {  
    If (TID [ c_val ].get_declare ( ))  
    st_lex.push (TID [ c_val ].get_type ( ) );  
    else throw "not declared";  
}
```

Для контролю контекстних умов кожної трійки - "операнд-операція-операнд" використовується функція `check_op`:

```
void Parser::check_op ( ) {
type_of_lex t1, t2, op, t = LEX_INT, r = LEX_BOOL;
t2 = st_lex.pop ( );
op = st_lex.pop ( );
t1 = st_lex.pop ( );
if (op == LEX_PLUS | op == LEX_MINUS || op == LEX_TIMES | == LEX_SLASH)
r = LEX_INT;
if (op == LEX_OR || op == LEX_AND)
t = LEX_BOOL;
if (t1 == t2 && t1 == t) st_lex.push( r );
else throw "wrong types are in operation";
prog.put_lex (Lex (op));
}
```

Для контролю за типом операнда односторонньої операції *not* будемо використовувати функцію `check_not`:

```
parser ::check_not ( ) {
if (st_lex.pop ( ) != LEX_BOOL)
throw "wrong type is in not";
else {
st_lex.push (LEX_BOOL);
```

Порівняємо граматики, що описують вирази, що складаються із символів + , \* , ( , ) , i:

G1:  $E \rightarrow E+E \mid E * E \mid (E) \mid i$  G4:  $E \rightarrow T \mid E+T$

$T \rightarrow F \mid T * F$

G2:  $E \rightarrow E+T \mid E * T \mid T F \rightarrow i \mid (E)$

$T \rightarrow i \mid (E)$

G3:  $E \rightarrow T+E \mid T * E \mid T$  G5:  $E \rightarrow T \mid T+E$

$T \rightarrow i \mid (E) \mid T \rightarrow F \mid F * T$

$F \rightarrow i \mid (E)$

# Правила виведення виразів модельної мови з діями для контролю контекстних умов

$$E \rightarrow E1 \mid E1 [= | < | >] \langle \text{st\_char.push}(\text{TD}[\text{c\_val}]) \rangle E1 \langle \text{check\_op}() \rangle$$
$$E1 \rightarrow T \{ [+ | - | \text{or} ] \langle \text{st\_char.push}(\text{TD}[\text{c\_val}]) \rangle T \langle \text{check\_op}() \rangle \}$$
$$T \rightarrow F \{ [* | / | \text{and} ] \langle \text{st\_char.push}(\text{TD}[\text{c\_val}]) \rangle F \langle \text{check\_op}() \rangle \}$$
$$F \rightarrow I \langle \text{check\_id}() \rangle \mid N \langle \text{st\_char.push}(\text{"int"}) \rangle \mid$$
$$[\text{true} \mid \text{false}] \langle \text{st\_char.push}(\text{"bool"}) \rangle \mid \text{not } F \langle \text{check\_not}() \rangle \mid (E)$$

## Контроль контекстних умов в операторах

$S \rightarrow I := E \mid \text{if } E \text{ then } S \text{ else } S \mid \text{while } E \text{ do } S \mid B \mid \text{read } (I) \mid \text{write } (E)$

### Оператор присвоєння $I:=E$

Контекстна умова: в операторі присвоєння типи змінної  $I$  та виразу  $E$  повинні збігатися.

В результаті контролю контекстних умов виразу  $E$  в стеку залишиться тип цього виразу (як тип результату останньої операції).

При аналізі ідентифікатора  $I$  перевіряється, чи він описаний, і його тип заноситься в той же стек (з допомогою функції `check_id()`). При цьому достатньо в потрібний момент рахувати з стека два елементи та порівняти їх:

```
parser::eq_type () {
if ( st_lex.pop() != st_lex.pop() )
throw "wrong types are in :=";
}
```

Правило виведення для оператора присвоєння:

$I \langle \text{check\_id} ( ) \rangle := E \langle \text{eq\_type} ( ) \rangle$

## *Умовний оператор, оператор циклу, оператор введення*

**if E then S else S | while E do S | read (I)**

### **Контекстні умови:**

- в умовному операторі та в операторі циклу як умова можливі лише логічні вирази.
- Операнд оператора введення повинен бути описаний.

Для контролю КУ в умовному операторі та операторі циклу функція `eq_bool ()`:

```
void Parser::eq_bool () {
if ( st_lex.pop() != LEX_BOOL )
throw "expression is not boolean";
}
```

Для перевірки операнда оператора введення `read (I)` використовується така функція:

```
void Parser::check_id_in_read () {
if ( !TID [c_val].get_declare( ) ) throw "not declared";
}
```

Правила виведення для умовного оператора та операторів циклу та введення:

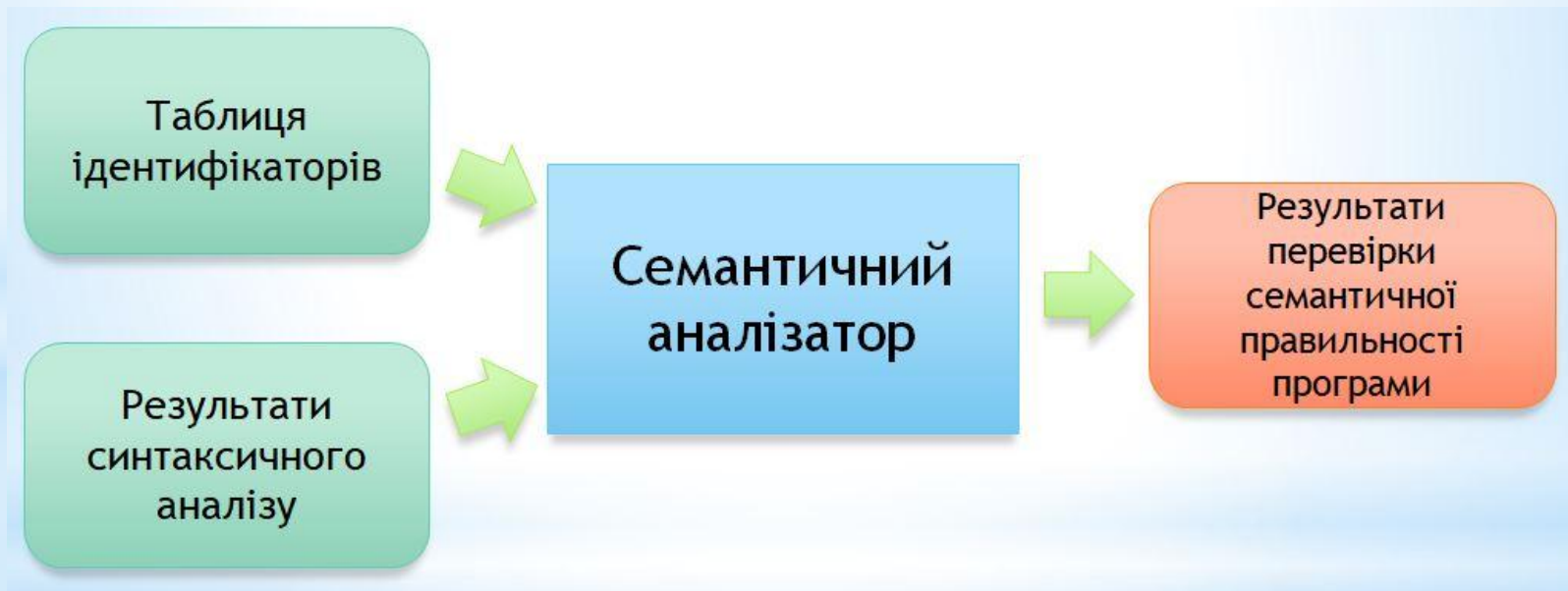
## Граматика з діями для розділу описів М-мови

```
void Parser::D () {
st_int.reset ();
if (c_type != LEX_ID) throw curr_lex;
else {
st_int.push (c_val);
gl ();
while (c_type == LEX_COMMA) {
gl ();
if (c_type != LEX_ID) throw curr_lex;
else {
st_int.push (c_val); gl ();
}
}
if (c_type != LEX_COLON) throw curr_lex;
else { gl ();
if (c_type == LEX_INT) { dec (LEX_INT); gl (); }
else
if (c_type == LEX_BOOL) { dec (LEX_BOOL); gl (); }
else throw curr_lex;
}
}
}
```

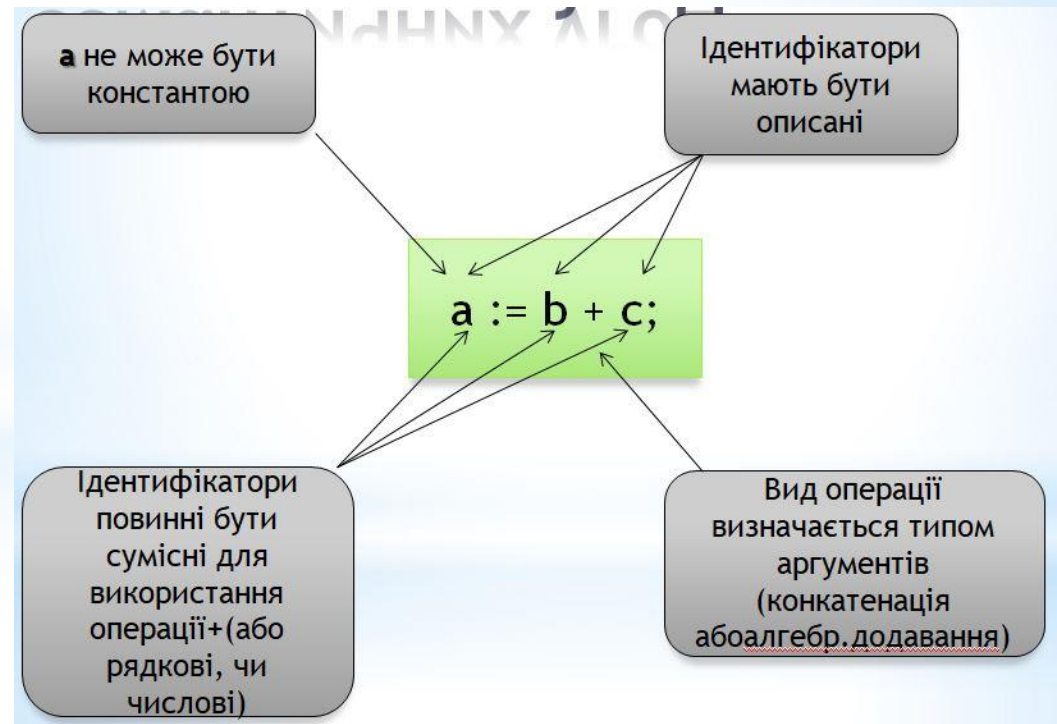
## 2. Семантичний аналізатор

**Семантичний аналізатор** виконує такі основні дії:

- перевірку дотримання у вхідній програмі семантичних угодами
- доповнення внутрішнього подання програми у компіляторі операторами та діями, неявно передбаченими семантикою вхідного мови
- перевірку елементарних семантичних (сміслових) норм мов програмування, які безпосередньо не пов'язані з



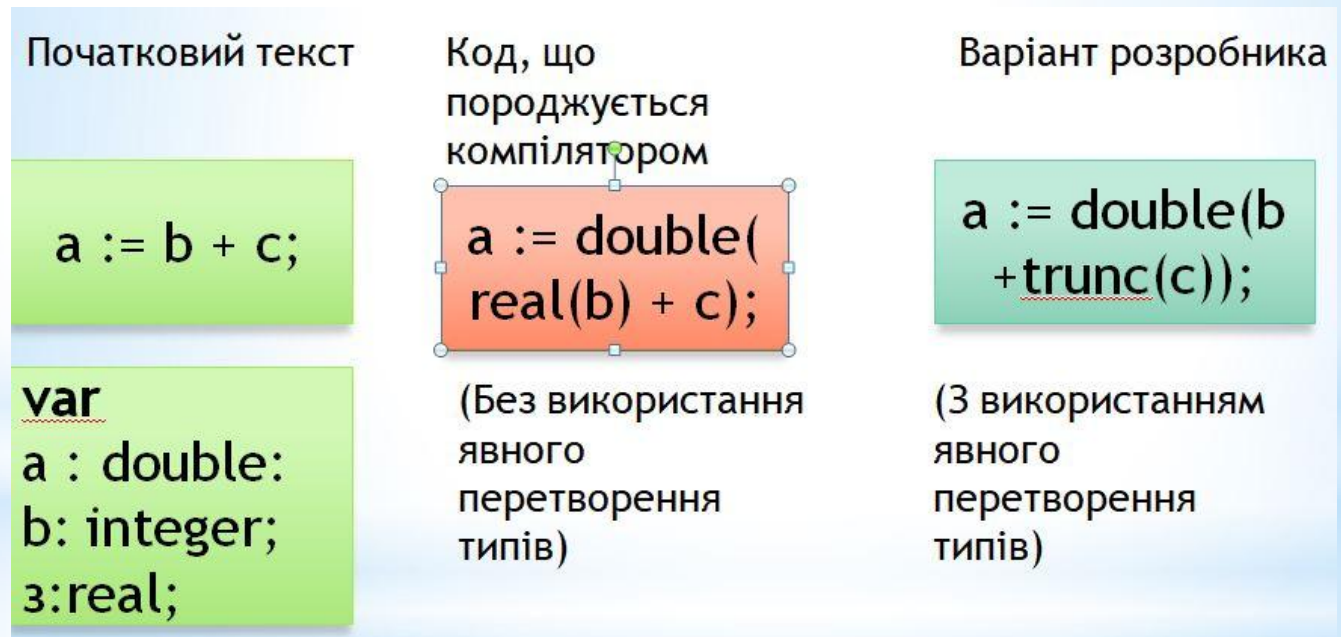
- кожна мітка, на яку є посилання, має один раз бути присутньою в програмі
- ідентифікатор повинен бути описаний один раз і жоден ідентифікатор не може бути описаний більш ніж один раз (з урахуванням блокової структури описів)
- всі операнди у виразах та операціях повинні мати типи, допустимі для даного виразу або операції
- типи змінних у виразах повинні бути узгоджені між собою
- при виклику процедур і функцій число та типи фактичних параметрів мають бути узгоджені з числом та типами формальних параметрів



*Перевірка дотримання семантичних угод*

Зв'язано з додаванням текст програми операторів та дій, неявно передбачених семантикою вхідної мови:

- перетворення типів операндів у виразах та при передачі параметрів у процедури та функції
- операції обчислення адреси, коли відбувається звернення до елементів складних структур даних



Доповнення внутрішнього представлення програми

## *Перевірка смислових норм мов програмування*

компілятором угод, виконання яких пов'язане зі змістом як усієї вихідної програми в цілому, так і окремих її фрагментів:

- кожна змінна або константа повинна хоча б один раз використовуватись у програмі;
- кожна змінна має бути визначена до її першого використання за будь-якого ходу виконання програми (першому використанню змінної має завжди передувати присвоєння їй будь-якого значення);
- результат функції має бути визначений за будь-якого ходу її виконання
- кожен оператор у вихідній програмі повинен мати можливість хоча б один раз виконатись;
- оператори умови та вибору повинні передбачати можливість ходу виконання програми з кожної зі своїх гілок;
- оператори циклу мають передбачати можливість завершення циклу
- ...

**Ідентифікація** змінних, типів, процедур, функцій та ін. лексичних одиниць мов програмування - це встановлення однозначної відповідності між лексичними одиницями та їх іменами у тексті вихідної програми

імена лексичних одиниць не повинні збігатися як між собою, так і з ключовими словами синтаксичних конструкцій мови  
локальні змінні мають область видимості

```
int f_test(int a)
{
    int b, c;
    b = 0;
    c = 0;
    if(b = 1){ return a;}
    c = a+b;
}
```

На етапі семантичного аналізу кожної лексичної одиниці мови видається унікальне ім'я в межах усієї вихідної програми і потім використовується при синтезі результуючої програми

- імена локальних змінних доповнюються іменами тих блоків (функцій, процедур), у яких ці змінні описані;

- імена внутрішніх змінних та функцій модулів вихідної програми доповнюються іменами самих модулів;

- імена процедур та функцій, що належать об'єктам (класам) в об'єктно-орієнтованих мовах програмування, доповнюються найменуваннями типів об'єктів (класів), яким вони належать;

- імена процедур та функцій модифікуються в залежності від типів їх формальних аргументів та ін.