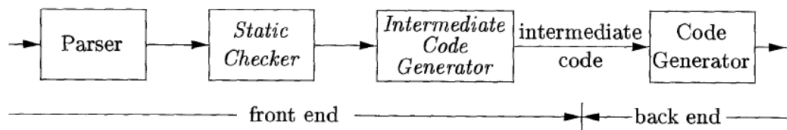


Intermediate-Code Generation

Intermediate-Code Generation



Three-Address Code

- In three-address code, there is at **most one operator** on the **right side of an instruction**
- Thus a **source-language expression** like $x+y*z$ might be translated into the sequence of three-address instructions

$$\begin{aligned}t_1 &= y * z \\t_2 &= x + t_1\end{aligned}$$

Common three-address instructions

1. Assignment instructions of the form $x = y \text{ op } z$, where op is a binary arithmetic or logical operation, and x , y , and z are addresses.
2. Assignments of the form $x = op \ y$, where op is a unary operation. Essential unary operations include unary minus, logical negation, shift operators, and conversion operators that, for example, convert an integer to a floating-point number.
3. *Copy instructions* of the form $x = y$, where x is assigned the value of y .
4. An unconditional jump `goto L` . The three-address instruction with label L is the next to be executed.
5. Conditional jumps of the form `if x goto L` and `ifFalse x goto L` . These instructions execute the instruction with label L next if x is true and false, respectively. Otherwise, the following three-address instruction in sequence is executed next, as usual.

6. Conditional jumps such as `if x relop y goto L` , which apply a relational operator (`<`, `==`, `>=`, etc.) to x and y , and execute the instruction with label L next if x stands in relation *relop* to y . If not, the three-address instruction following `if x relop y goto L` is executed next, in sequence.
7. Procedure calls and returns are implemented using the following instructions: `param x` for parameters; `call p, n` and `y = call p, n` for procedure and function calls, respectively; and `return y` , where y , representing a returned value, is optional. Their typical use is as the sequence of three-address instructions

```
param  $x_1$ 
param  $x_2$ 
...
param  $x_n$ 
call  $p, n$ 
```

generated as part of a call of the procedure $p(x_1, x_2, \dots, x_n)$. The integer n , indicating the number of actual parameters in “`call p, n` ,”

8. Indexed copy instructions of the form $x = y[i]$ and $x[i] = y$. The instruction $x = y[i]$ sets x to the value in the location i memory units beyond location y . The instruction $x[i] = y$ sets the contents of the location i units beyond x to the value of y .

Common three-address instructions

```
do i = i+1; while (a[i] < v);
```

```
L:  t1 = i + 1  
    i = t1  
    t2 = i * 8  
    t3 = a [ t2 ]  
    if t3 < v goto L
```

Data structures for TAC

Quadruples

A *quadruple* (or just “*quad*”) has four fields, which we call *op*, *arg*₁, *arg*₂, and *result*. The *op* field contains an internal code for the operator. For instance, the three-address instruction $x = y + z$ is represented by placing $+$ in *op*, *y* in *arg*₁, *z* in *arg*₂, and *x* in *result*. The following are some exceptions to this rule:

1. Instructions with unary operators like $x = \text{minus } y$ or $x = y$ do not use *arg*₂. Note that for a copy statement like $x = y$, *op* is $=$, while for most other operations, the assignment operator is implied.
2. Operators like *param* use neither *arg*₂ nor *result*.
3. Conditional and unconditional jumps put the target label in *result*.

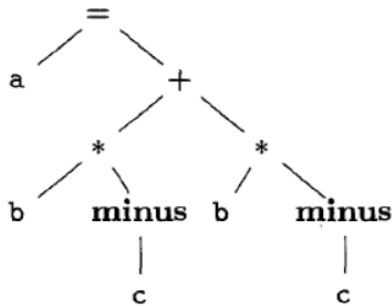
Data structures for TAC

Quadruples

Three-address code for the assignment $a = b * -c + b * -c$;

```
t1 = minus c
t2 = b * t1
t3 = minus c
t4 = b * t3
t5 = t2 + t4
a = t5
```

(a) Three-address code



(a) Syntax tree

Data structures for TAC

Quadruples

Three-address code for the assignment $a = b * -c + b * -c$;

```
t1 = minus c
t2 = b * t1
t3 = minus c
t4 = b * t3
t5 = t2 + t4
a = t5
```

	<i>op</i>	<i>arg₁</i>	<i>arg₂</i>	<i>result</i>
0	minus	c		t ₁
1	*	b	t ₁	t ₂
2	minus	c		t ₃
3	*	b	t ₃	t ₄
4	+	t ₂	t ₄	t ₅
5	=	t ₅		a
	...			

(a) Three-address code

(b) Quadruples

Data structures for TAC

Triples

- A triple has **only three fields**, which we call **op**, **arg1**, and **arg2**
- Note that the **result field in Quad** is used primarily for temporary names.
- Using **triples**, we refer to the **result of an operation** $x \text{ op } y$ **by its position**, rather than by an explicit temporary name.
- Thus, instead of the **temporary t**, a triple representation would refer to **position (0)**.
- **Parenthesized numbers** represent pointers into the **triple structure** itself.

Data structures for TAC

Triples

Three-address code for the assignment $a = b * -c + b * -c$;

```
t1 = minus c
t2 = b * t1
t3 = minus c
t4 = b * t3
t5 = t2 + t4
a = t5
```

(a) Three-address code

	<i>op</i>	<i>arg₁</i>	<i>arg₂</i>
0	minus	c	
1	*	b	(0)
2	minus	c	
3	*	b	(2)
4	+	(1)	(3)
5	=	a	(4)
	...		

(b) Triples

Benefit of Quad over Triples

- A benefit of **quadruples** over triples can be seen in an **optimizing compiler**, where **instructions are often moved around**.
 - With quadruples, **if we move an instruction** that **computes a temporary t**, then the instructions that **use t** require **no change**.
- With **triples**, the **result** of an operation is referred to by **its position**
- So **moving an instruction** may require us to **change all references** to that result.

Indirect triples

- Consist of a **listing of pointers** to triples,
 - Rather than a listing of triples themselves.
- For example, use an **array *instruction*** to list **pointers to triples** in the desired order.

With indirect triples, an **optimizing compiler** can move an instruction by **reordering the instruction list**, without affecting the triples themselves.

<i>instruction</i>		<i>op</i>	<i>arg₁</i>	<i>arg₂</i>
35	(0)	minus	c	
36	(1)	*	b	(0)
37	(2)	minus	c	
38	(3)	*	b	(2)
39	(4)	+	(1)	(3)
40	(5)	=	a	(4)
	...			

0	minus	c	
1	*	b	(0)
2	minus	c	
3	*	b	(2)
4	+	(1)	(3)
5	=	a	(4)
		...	

Static Single-Assignment Form

Two distinctive aspects distinguish SSA from three-address code.

(a) The first is that **all assignments in SSA** are to **variables with distinct names**; hence the term static single-assignment.

$p = a + b$

$q = p - c$

$p = q * d$

$p = e - p$

$q = p + q$

$p_1 = a + b$

$q_1 = p_1 - c$

$p_2 = q_1 * d$

$p_3 = e - p_2$

$q_2 = p_3 + q_1$

(a) Three-address code. (b) Static single-assignment form.

Static Single-Assignment Form

Two distinctive aspects distinguish SSA from three-address code.

(a) The first is that **all assignments in SSA** are to **variables with distinct names**; hence the term static single-assignment.

(b)

The same variable may be defined in two different control-flow paths in a program. For example, the source program

```
if ( flag ) x = -1; else x = 1;  
y = x * a;
```

has two control-flow paths in which the variable x gets defined. If we use different names for x in the true part and the false part of the conditional statement, then which name should we use in the assignment $y = x * a$? Here is where the second distinctive aspect of SSA comes into play. SSA uses a notational convention called the ϕ -function to combine the two definitions of x :

```
if ( flag )  $x_1 = -1$ ; else  $x_2 = 1$ ;  
 $x_3 = \phi(x_1, x_2)$ ;
```

Static Single-Assignment Form

Here, $\phi(x_1, x_2)$ has the value x_1 if the control flow passes through the true part of the conditional and the value x_2 if the control flow passes through the false part. That is to say, the ϕ -function returns the value of its argument that corresponds to the control-flow path that was taken to get to the assignment-statement containing the ϕ -function.

Major translation classes of Three address code generation

- (a) Declaration statements (+ handling data type and storage)
- (b) Expressions and statements
- (a) Control flow statements

Declaration statement

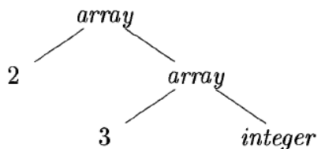
Representing data types: **Type Expressions**

Types have structure, which we shall represent using type expressions.

- A **type expression** is either a **basic type** (*boolean, char, integer, float, and void*)
or
- is formed by **applying an operator** called a **type constructor** to a type expression.
- A **type expression** can be formed by applying the **array type constructor** to a **number** and a **type expression**.

Declaration statement

- The array type `int [2] [3]` can be read as "array of 2 arrays of 3 integers each"
- Can be represented as a **type expression** `array(2, array(3, integer))`.
- This type is represented by the tree.



- The **operator array** takes **two parameters**, a number and a type.
 - Here the **type expression** can be formed by applying the **array type constructor** to a number and a type expression.

Declaration statement

Example SDD

PRODUCTION	SEMANTIC RULES
$T \rightarrow B C$	$T.t = C.t$ $C.b = B.t$
$B \rightarrow \text{int}$	$B.t = \text{integer}$
$B \rightarrow \text{float}$	$B.t = \text{float}$
$C \rightarrow [\text{num}] C_1$	$C.t = \text{array}(\text{num.val}, C_1.t)$ $C_1.b = C.b$
$C \rightarrow \epsilon$	$C.t = C.b$



Type Expressions

- Nonterminal T generates either a **basic type** or an **array type**.
- Nonterminal B generates one of the basic types int and float.
- T generates a basic type when C derives ϵ .
- Otherwise, C generates array components consisting of a sequence of integers, each integer surrounded by brackets.

Declaration statement

Example SDD

PRODUCTION	SEMANTIC RULES
$T \rightarrow B C$	$T.t = C.t$ $C.b = B.t$
$B \rightarrow \text{int}$	$B.t = \text{integer}$
$B \rightarrow \text{float}$	$B.t = \text{float}$
$C \rightarrow [\text{num}] C_1$	$C.t = \text{array}(\text{num.val}, C_1.t)$ $C_1.b = C.b$
$C \rightarrow \epsilon$	$C.t = C.b$

← **Type Expressions**

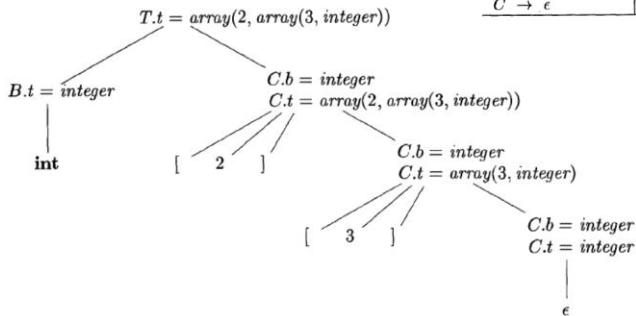
- The nonterminals B and T have a synthesized attribute t representing a type.
- The nonterminal C has two attributes: an inherited attribute b and a synthesized attribute t .

Declaration statement

Example SDD

input string `int [2][3]`

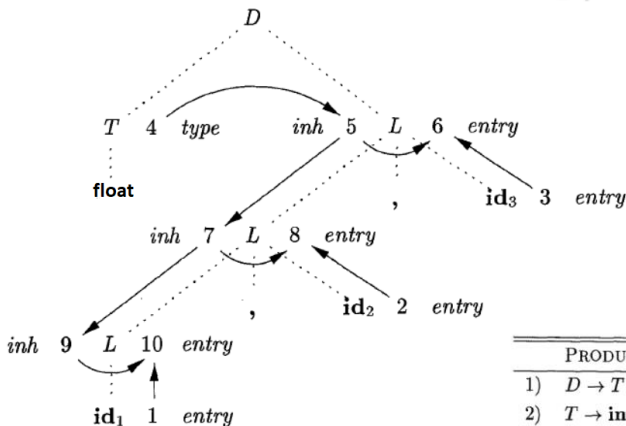
PRODUCTION	SEMANTIC RULES
$T \rightarrow B C$	$T.t = C.t$ $C.b = B.t$
$B \rightarrow \text{int}$	$B.t = \text{integer}$
$B \rightarrow \text{float}$	$B.t = \text{float}$
$C \rightarrow [\text{ num }] C_1$	$C.t = \text{array}(\text{num.val}, C_1.t)$ $C_1.b = C.b$
$C \rightarrow \epsilon$	$C.t = C.b$



- The nonterminal C has two attributes: an inherited attribute b and a synthesized attribute t .
- The inherited b attributes pass a basic type down the tree, and the synthesized t attributes accumulate the result.

Declaration statement: Example SDD

float id_1 , id_2 , id_3



PRODUCTION	SEMANTIC RULES
1) $D \rightarrow T L$	$L.inh = T.type$
2) $T \rightarrow \text{int}$	$T.type = \text{integer}$
3) $T \rightarrow \text{float}$	$T.type = \text{float}$
4) $L \rightarrow L_1 , \text{id}$	$L_1.inh = L.inh$ $addType(id.entry, L.inh)$
5) $L \rightarrow \text{id}$	$addType(id.entry, L.inh)$

Symbol table

ST(global)

This is the Symbol Table for global symbols

Name	Type	Initial Value	Size	Offset	Nested Table
d	float	2.3	8	0	null
i	int	null	4	8	null
w	<i>array</i> (10, int)	null	40	12	null
a	int	4	4	52	null
p	<i>ptr</i> (int)	null	4	56	null
b	int	null	4	60	null
func	<i>function</i>	null	0	64	ptr-to-ST(func)
c	char	null	1	64	null

Find the storage for each variable

More on Declaration statement

Data type + Storage layout

$$\begin{aligned} D &\rightarrow T \text{ id} ; D \mid \epsilon \\ T &\rightarrow B C \\ B &\rightarrow \text{int} \mid \text{float} \\ C &\rightarrow \epsilon \mid [\text{num}] C \end{aligned}$$

- Simplified grammar that declares just **one name at a time**;
- We already explored the declarations with lists of names

Storage layout:

- **Relative address** of all the variables
- From the **type of a name**, we can determine the **amount of storage** that will be needed for the name at run time.
- At **compile time**, we can use these amounts to **assign each name a relative address**.
- The **type** and **relative address** are saved in the **symbol-table entry** for the name.

More on Declaration statement

Data type + Storage layout

- The **width** of a **type** is the number of **storage units** needed for objects of that type (**offset**).
- A **basic type**, such as a character, integer, or float, requires an integral number of bytes.
- **Arrays** allocated in one **contiguous block of bytes**

$$T \rightarrow \begin{matrix} B \\ C \end{matrix} \quad \{ t = B.type; w = B.width; \}$$
$$B \rightarrow \text{int} \quad \{ B.type = \text{integer}; B.width = 4; \}$$
$$B \rightarrow \text{float} \quad \{ B.type = \text{float}; B.width = 8; \}$$
$$C \rightarrow \epsilon \quad \{ C.type = t; C.width = w; \}$$
$$C \rightarrow [\text{num}] C_1 \quad \{ \text{array}(\text{num.value}, C_1.type); \\ C.width = \text{num.value} \times C_1.width; \}$$

Computes **data types** and their **widths** for basic and array types

 **Type Expressions**

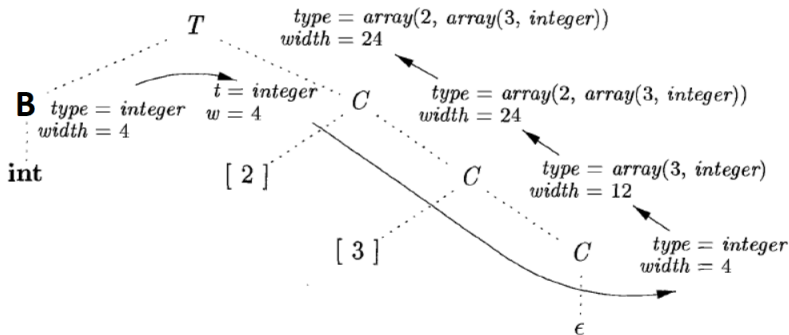
The **width of an array** is obtained by **multiplying** the **width of an element** by the **number of elements** in the array.

More on Declaration statement

Data type + Storage layout

`int[2][3]`

$T \rightarrow B$ C	$\{ t = B.type; w = B.width; \}$
$B \rightarrow \text{int}$	$\{ B.type = \text{integer}; B.width = 4; \}$
$B \rightarrow \text{float}$	$\{ B.type = \text{float}; B.width = 8; \}$
$C \rightarrow \epsilon$	$\{ C.type = t; C.width = w; \}$
$C \rightarrow [\text{num}] C_1$	$\{ \text{array}(\text{num.value}, C_1.type);$ $C.width = \text{num.value} \times C_1.width; \}$



Relative address

Name	Data type
d	float
i	int
w	<i>array</i> (10, int)

Relative address

0
8
12
16

Sequences of Declarations

```
int x;  
float y;
```

- Relative address: **offset**
- Keeps track of the **next available relative address**

$$P \rightarrow \{ \text{offset} = 0; \} D$$
$$D \rightarrow T \text{ id} ; \quad \{ \text{top.put}(\text{id.lexeme}, T.\text{type}, \text{offset}); \\ \text{offset} = \text{offset} + T.\text{width}; \}$$
$$D \xrightarrow{D_1} \epsilon$$

- The translation scheme deals with a **sequence of declarations** of the form $T \text{ id}$, where T generates a data type
- Before the first declaration is considered, **offset is set to 0**.
- As **each new name x** is seen, x is entered into the **symbol table** with its **relative address = current value of offset**,
 - which is **then incremented** by the width of the type of x .

Sequences of Declarations

- Relative address: offset
- Keeps track of the next available relative address

$$P \rightarrow \{ \text{offset} = 0; \} D$$
$$D \rightarrow T \text{ id} ; \quad \{ \text{top.put}(\text{id.lexeme}, T.type, \text{offset}); \\ \text{offset} = \text{offset} + T.width; \}$$
$$D \xrightarrow{D_1} \epsilon$$

The semantic action within the production $D \rightarrow T \text{ id} ; D_1$ creates a symbol-table entry by executing $\text{top.put}(\text{id.lexeme}, T.type, \text{offset})$. Here top denotes the current symbol table. The method top.put creates a symbol-table entry for id.lexeme , with type $T.type$ and relative address offset in its data area.

Major translation classes of Three address code generation

(a) Declaration statements (+ handling data type and storage)

(b) Expressions and statements

(a) Control flow statements

Translation of Expressions

statement $a = b + -c$

Three-address code for an assignment statement

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.\text{code} = E.\text{code} $ $\text{gen}(\text{top.get}(\text{id.lexeme}) '=' E.\text{addr})$
$E \rightarrow E_1 + E_2$	$E.\text{addr} = \text{new Temp}()$ $E.\text{code} = E_1.\text{code} E_2.\text{code} $ $\text{gen}(E.\text{addr} '=' E_1.\text{addr} '+' E_2.\text{addr})$
$ - E_1$	$E.\text{addr} = \text{new Temp}()$ $E.\text{code} = E_1.\text{code} $ $\text{gen}(E.\text{addr} '=' \text{'minus'} E_1.\text{addr})$
$ (E_1)$	$E.\text{addr} = E_1.\text{addr}$ $E.\text{code} = E_1.\text{code}$
$ \text{id}$	$E.\text{addr} = \text{top.get}(\text{id.lexeme})$ $E.\text{code} = ''$

Attribute **code** for S

attributes **addr** and **code** for an expression E.

Attributes **S.code** and **E.code** denote the **three-address code** for S and E, respectively.

Attribute **E.addr** denotes the **address** that will hold the **value of E**.

Translation of Expressions

Three-address code for an assignment statement

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.\text{code} = E.\text{code} $ $\text{gen}(\text{top.get}(\text{id.lexeme}) '=' E.\text{addr})$
$E \rightarrow E_1 + E_2$	$E.\text{addr} = \text{new Temp}()$ $E.\text{code} = E_1.\text{code} E_2.\text{code} $ $\text{gen}(E.\text{addr} '=' E_1.\text{addr} '+' E_2.\text{addr})$
$ - E_1$	$E.\text{addr} = \text{new Temp}()$ $E.\text{code} = E_1.\text{code} $ $\text{gen}(E.\text{addr} '=' \text{'minus'} E_1.\text{addr})$
$ (E_1)$	$E.\text{addr} = E_1.\text{addr}$ $E.\text{code} = E_1.\text{code}$
$ \text{id}$	$E.\text{addr} = \text{top.get}(\text{id.lexeme})$ $E.\text{code} = ''$

When an expression is a **single identifier**, say **x**, then **x itself holds the value of the expression**.

The semantic rules for this production define **E.addr** to point to the symbol-table entry

Translation of Expressions

Three-address code for an assignment statement

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.code = E.code \parallel$ $gen(top.get(id.lexeme) '=' E.addr)$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $gen(E.addr '=' E_1.addr '+' E_2.addr)$



The semantic rules for $E \rightarrow E_1 + E_2$, generate code to compute the value of E from the values of E_1 and E_2 . Values are computed into newly generated temporary names. If E_1 is computed into $E_1.addr$ and E_2 into $E_2.addr$, then $E_1 + E_2$ translates into $t = E_1.addr + E_2.addr$, where t is a new temporary name. $E.addr$ is set to t . A sequence of distinct temporary names t_1, t_2, \dots is created by successively executing **new Temp()**.

Translation of Expressions

Three-address code for an assignment statement

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.code = E.code \parallel$ $gen(top.get(\text{id.lexeme}) '=' E.addr)$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $gen(E.addr '=' E_1.addr '+' E_2.addr)$

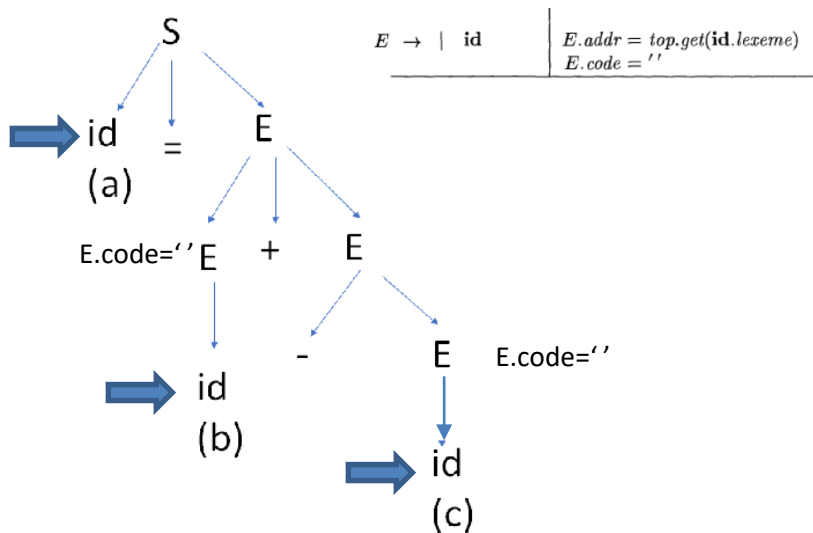
Finally, the production $S \rightarrow \text{id} = E ;$ generates instructions that assign the value of expression E to the identifier id . The semantic rule for this production uses function $top.get$ to determine the address of the identifier represented by id , as in the rules for $E \rightarrow \text{id}$. $S.code$ consists of the instructions to compute the value of E into an address given by $E.addr$, followed by an assignment to the address $top.get(\text{id.lexeme})$ for this instance of id .

Translation of Expressions

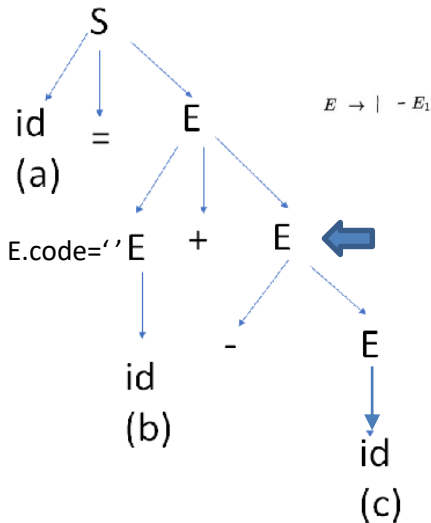
Three-address code for an assignment statement

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.code = E.code \parallel$ $gen(top.get(\text{id.lexeme}) '=' E.addr)$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $gen(E.addr '=' E_1.addr '+' E_2.addr)$
$ - E_1$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel$ $gen(E.addr '=' \text{'minus'} E_1.addr)$
$ (E_1)$	$E.addr = E_1.addr$ $E.code = E_1.code$
$ \text{id}$	$E.addr = top.get(\text{id.lexeme})$ $E.code = ''$

statement $a = b + - c$



statement `a = b + - c`



$$E \rightarrow | - E_1$$

$$E.addr = \text{new Temp}()$$

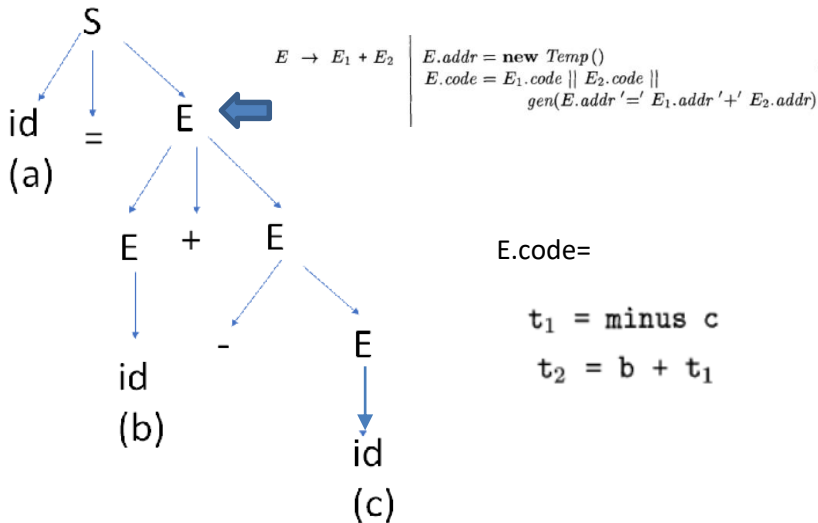
$$E.code = E_1.code ||$$

$$gen(E.addr = 'minus' E_1.addr)$$

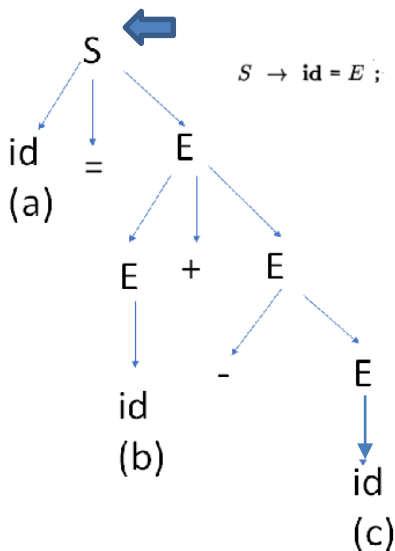
`E.code=`

`t1 = minus c`

statement `a = b + - c`



statement $a = b + - c$



$$S \rightarrow \mathbf{id} = E ; \quad \left\{ \begin{array}{l} S.code = E.code \parallel \\ gen(top.get(\mathbf{id.lexeme}) \neq E.addr) \end{array} \right.$$

S.code=

```
t1 = minus c
t2 = b + t1
a = t2
```

statement `a = b + - c`

Translation of Expressions

Three-address code for an assignment statement

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.code = E.code \parallel$ $gen(top.get(\text{id.lexeme}) '=' E.addr)$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $gen(E.addr '=' E_1.addr '+' E_2.addr)$
$ - E_1$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel$ $gen(E.addr '=' \text{'minus'} E_1.addr)$
$ (E_1)$	$E.addr = E_1.addr$ $E.code = E_1.code$
$ \text{id}$	$E.addr = top.get(\text{id.lexeme})$ $E.code = ''$

statement $a = b + - c$

$t_1 = \text{minus } c$

$t_2 = b + t_1$

$a = t_2$

Incremental Translation

- So far, ***E.Code*** attributes were long strings
 - Generated incrementally
- In incremental translation, generate only the **new three-address instructions**
- **Past sequence** may either be **retained in memory** for further processing, or it may be **output incrementally**.
- In the incremental approach, `gen()` not only constructs a three-address instruction,
 - it **appends** the instruction to the sequence of instructions generated so far.

Incremental Translation

$$S \rightarrow \text{id} = E ; \quad \{ \text{gen}(\text{top.get}(\text{id.lexeme}) \text{'=' } E.\text{addr}); \}$$
$$E \rightarrow E_1 + E_2 \quad \{ E.\text{addr} = \text{new Temp}(); \\ \text{gen}(E.\text{addr} \text{'=' } E_1.\text{addr} \text{'+' } E_2.\text{addr}); \}$$
$$\mid - E_1 \quad \{ E.\text{addr} = \text{new Temp}(); \\ \text{gen}(E.\text{addr} \text{'=' } \text{'minus' } E_1.\text{addr}); \}$$
$$\mid (E_1) \quad \{ E.\text{addr} = E_1.\text{addr}; \}$$
$$\mid \text{id} \quad \{ E.\text{addr} = \text{top.get}(\text{id.lexeme}); \}$$

- This **translation scheme generates the same code** as the previous syntax directed definition.
- With the incremental approach, the **E.code attribute is not used**,
 - Since there is a **single sequence of instructions** that is created by **successive calls to gen()**.