

# Apuntes de Compiladores

Mauricio Elían Delgadillo García



# Índice general

|   |           |
|---|-----------|
| <b>1. Conceptos Introductorios</b>                              | <b>5</b>  |
| 1.1. 1er Paso <u>El Assembler</u> (ASM)                         | 5         |
| <b>2. Lenguaje Intermedio</b>                                   | <b>7</b>  |
| 2.1. Porque usar I.L.?  | 7         |
| 2.2. ElCodigo de 3 direcciones                                  | 8         |
| 2.2.1. Operaciones Arimeticas y Logicas                         | 8         |
| 2.2.2. Usando operaciones relacionales (OPREL)                  | 10        |
| 2.2.3. Incremento - Decrementos                                 | 10        |
| 2.2.4. Saltos (Jump's or GOTO's)                                | 10        |
| 2.2.5. Operaciones del System (S.O.)                            | 11        |
| 2.3. Procedimientos en el Codigo-3                              | 12        |
| 2.3.1. Las instrucciones CALL y RET                             | 12        |
| 2.3.2. Las variables temporales ("Auxiliares" del compiladores) | 12        |
| 2.4. Esquema de Traducción                                      | 14        |
| 2.4.1. Las expresiones (booleanas y aritméticas)                | 14        |
| 2.4.2. La Asignación  | 15        |
| 2.4.3. Traducción para el If y While                            | 15        |
| 2.5. Representación del Codigo-3                                | 16        |
| 2.6. Tabla de Simbolos  | 18        |
| 2.6.1. La T.S.I.D.  | 18        |
| 2.6.2. La T.S.S.  | 19        |
| 2.6.3. Codificaciones de las direcciones                        | 19        |
| <b>3. Analizador Léxico</b>                                     | <b>21</b> |
| 3.1. Manejo del Programa Fuente                                 | 22        |
| 3.1.1. El fin de linea (EOLN)                                   | 24        |
| 3.2. Especificación de los Componentes Lexicos                  | 25        |
| 3.2.1. Definiciones Regulares                                   | 26        |
| 3.2.2. Lexemas (Lexem's)  | 27        |
| 3.2.3. Números Reales   | 27        |
| 3.3. Diagrama de Transiciones (dt)                              | 28        |
| 3.3.1. El ANALEX como un jugador de sopa de letras              | 28        |
| 3.3.2. El Espacio   | 30        |
| 3.3.3. Los Tokens ERROR y FIN                                   | 30        |

|                                   |    |
|-----------------------------------|----|
| 3.4. “Reglas” de los dt . . . . . | 31 |
|-----------------------------------|----|

# Capítulo 1

## Conceptos Introductorios

Como se sabe, la computadora solo entiende 1's y 0's (Lenguaje Maquina)

(IMAGEN)

En el principio, el programador debía escribir su código (Programa) en Binario.

### 1.1. 1er Paso El Assembler (ASM)

Lo que pasa en el Assembler es que dado una especie de alias dándole una instrucción.

En vez de escribir : 101011010  
Se escribiría : MOV EAX, EAX

La conversión del ASM a lenguaje Maquina (Binaria) es 1 a 1, osea equivalente.

1 Instrucción ASM = 1 Instrucción (del lenguaje) maquina

mientras un lenguaje esta mas cercano al lenguaje maquina, se dice **Lenguaje de Bajo Nivel**

(IMAGEN )

El primer lenguaje de programación de Alto Nivel fue FORTRAN (**Form**ula **Tr**anslation) el cual se tardo 18 años en completarse (El compilador)

(IMAGEN )

Un compilador de un lenguaje de Alto Nivel traduce una instrucción del lenguaje a N instrucciones maquina.

(IMAGEN )

**Las 3 instrucciones básicas:** IF - THEN, IF - THEN - ELSE, WHILE (y la Asignación  $x:=3$ ;) )

# Capítulo 2

## Lenguaje Intermedio

/\* Lenguaje Intermedio = Intermediate Language = IL \*/  
Recordemos que:

(IMAGEN )

El IL es un lenguaje de bajo nivel, muy cercano al Assembler, pero no es entendido por ninguna computadora del mundo.

### 2.1. Porque usar I.L.?

Un lenguaje de alto nivel, es muy complicado para llevarlo al lenguaje maquina. Por este motivo, el diseñador del compilador se inventa un lenguaje que le permite, mas fácilmente hacer la traducción:

(IMAGEN )

Otra razón de la asistencia del I.L., es la de poder generar archivos ejecutables en distintas **plataformas de Hw y Sw** (Arquitectura propia de la computadora, su lenguaje maquina y sistemas operativos)

/\* Archivo Ejecutable = que es capaz de correr solo (Stand-alone)

El archivo ejecutable (Aplicación o App) contiene código maquina e.g. Windows el .exe

(IMAGEN )

\*/

Ademas, que es la moda Actual, al I.L. se le crea un interprete que ejecute el lenguaje intermedio generado. A este Interprete se lo llama Virtual Machina o V.M.

//V.M. = Computadora Imaginaria

### Cual es el problema de los V.M.?

La V.M., obviamente es una APP (un ejecutable) escrito para la plataforma donde va a correr (e.g. Windows, Linux, MacOS)

Por ejemplo JAVA:

(IMAGEN)

En Windows instalo la V.M. de Java, llamado Java Virtual Machine (J.V.M.), que viene en el archivo Java.exe

(IMAGEN)

**Los IDE's no compilan nada ni corren ningún programa**, Para Linux, necesitamos el ejecutable Java

(IMAGEN)

## 2.2. El Codigo de 3 direcciones

Es el I.L. mas utilizado en la compilación. En realidad, todos los I.L.'s se basan en el.

Se llama de 3 direcciones, porque a lo sumo cada una de sus instrucciones ocupa 3 direcciones de memoria. Por direcciones de memoria entendemos a variables, procedimientos etiquetas

Ejemplo:

(IMAGEN)

### 2.2.1. Operaciones Arimeticas y Logicas

Toda operación aritmética o lógica, almacena su resultado en una variable:

$$\text{Var1} = \text{Var2} \text{ oper. } \text{Var3}$$

- Var1,Var2,Var3 = Variables
- oper. = operaciones =  $\{*, +, -, /, \text{mod}, \text{and}, \text{or}\}$

Ejemplo:



$$x = z - y$$

Las operaciones (cálculos) se hacen entre variables.

Por ejemplo:

$x = z + \underline{2} \rightarrow 2$  es una constante no puede ejecutar esta operación, debe ser una variable.

$$x = -y \text{ //operación minus}$$

También se permite la asignación simple:

$$p = 25 \text{ //asignación a una ctte.}$$

$$x = p \text{ //asignación a una variable.}$$

En el caso lógico

$\{\sim, \vee, \wedge\}$  se puede generar cualquier conectivo entre ellos.

$$\begin{aligned} p &= \text{not } q \\ q &= z \text{ and } x \\ w &= p \text{ or } y \end{aligned}$$

No existe tipo de datos, son solo números, no existe el **boolean** en si.

Como el C3 solo maneja números, hace un tratamiento lógico, usando el criterio del lenguaje C.

| Cuando lee                     | Cuando escribe |
|--------------------------------|----------------|
| Si es $!= 0 \rightarrow true$  | $1 = true$     |
| Si es $== 0 \rightarrow false$ | $0 = false$    |

Por ejemplo:

$$\begin{aligned} x &= 20 \\ y &= 0 \\ z &= x \text{ or } y \end{aligned}$$

Cuanto vale z?

$$\begin{aligned} x = 20 \neq 0 &= true \\ y = 0 &= false \\ z = true \text{ or } false &= true = 1 \\ \textbf{Respuesta: } z &= 1 \end{aligned}$$

### 2.2.2. Usando operaciones relacionales (OPREL)

// <, >, ≤, ≥, =, ≠

$x = (y \geq z) \rightarrow x = \text{valor } (0 \text{ o } 1) \text{ de que } y \text{ es mayor o igual que } z$

Por ejemplo:

$$\begin{aligned} t_1 &= 1 \\ t_2 &= 5 \\ z &= (t_1 = t_2) \text{ // } z = 0 \text{ (false)} \\ y &= (t_1 \geq t_2) \text{ // } y = 1 \text{ (true)} \end{aligned}$$

### 2.2.3. Incremento - Decrementos

$$\begin{array}{ll} x = 2 & \\ \text{INC } x \text{ //Equivalente a } x++ & \text{DEC } x \text{ //Equivalente a } x-- \\ t_1 = x \text{ // } t_1 = 3 & t_1 = x \text{ // } t_1 = 1 \end{array}$$

### 2.2.4. Saltos (Jump's or GOTO's)

**Saltos Incondicionales**

GOTO Etiqueta

Las etiquetas pueden ser :

$E_1 : E_2 : E_3 : \dots$

Las etiquetas no pueden tomar cualquier nombre

**Saltos Condicionales**

Un salto condicional se hace usando el valor de verdad de una variable. Entonces, Se tienen dos posibles caminos:

- if(var = 1)  $\Rightarrow$  GOTO etiqueta  
 //Si var es true (osea distinto de 0) ira a la etiqueta  
 //Si var es false (osea igual a 0) esta linea no se ejecuta, osea seguiria bajando
- if(var = 0)  $\Rightarrow$  GOTO etiqueta  
 //Si var es true (osea distinto de 0) esta linea no se ejecuta, osea seguiria bajando  
 //Si var es false (osea igual a 0) ira a la etiqueta

Ejemplo:

**Leer un numero N (Read(N)) y mostrar N asteriscos en la pantalla**

Introduzca N: **5**  
 \*\*\*\*\*

*Solución:*

En Delphi:

```
writes ("Introduzca _N:")
READ(N)
for i=1 to N Do{
  writes (*)
}
```

En C3:

```
writes ("Introduzca _N:")
READ(N)
i = 1
E1: t1 = (i ≤ N)
    if (t1 = 0) ⇒ Goto E2 //Va a E2 si t1 es False
    writes ("x")
    inc i
    Goto E1
E2:
```

### 2.2.5. Operaciones del System (S.O.)

- Para la Salida

|  |  |  |
|--|--|--|
| <pre>-Writes ("Mensaje") //1 solo mensaje //(String ctte.)</pre> | <pre>-Write(var) //1 sola variable</pre> | <pre>-NL //New Line //(para que el cursor //baje a la siguiente //linea)</pre> |
|--|--|--|

- Para la entrada

```
-Read(var)
//1 sola variable
```

Ejemplos:

| Codigo-3   | Resultado             |
|--|-----------------------|
| $\beta$ = paso en blanco<br>$x = 4$<br>Writes ("El_valor_de_x_es_ $\beta$ ")<br>Write(x)       | El valor de x es 4    |
| $\beta$ = paso en blanco<br>$x = 4$<br>Writes ("El_valor_de_x_es_ $\beta$ ")<br>NL<br>Write(x) | El valor de x es<br>4 |

- Leer un numero N y calcular  $F = N!$

| Codigo-3   | Resultado                |
|--|--------------------------|
| Read(N)<br>$F = 1$<br>$i = 1$<br>E1: $t1 = (i \leq N)$<br>$t2 = (t1 = 0) \Rightarrow \text{Goto E2}$<br>$F = F * i$<br>INC i<br>Goto E1<br>E2: Writes ("El_Factorial_de_ $\beta$ ")<br>Write(N)<br>Writes (" $\beta$ es $\beta$ ")<br>Write(F) | El Factorial de 5 es 120 |

## 2.3. Procedimientos en el Codigo-3

Como todo programa, el código de tres direcciones (Codigo-3, C3) se divide en procedimientos, teniendo al \$MAIN como el principal.

Es decir, un programa C-3 tiene como mínimo un procedimiento: \$MAIN

### 2.3.1. Las instrucciones CALL y RET

| Instrucción | Equivalente (Java) |
|-------------|--------------------|
| CALL Proc   | //Proc();          |
| RET         | // Return;         |

### 2.3.2. Las variables temporales (“Auxiliares” del compiladores)

Las variables temporales son “variables” creadas por el compilador para realizar calculos intermedios durante la traducción. Las variables temporales, también se crean en la pila de la aplicación y son de uso local.

| Programador                                       | Codigo - 3   |
|---|--|
| <pre> <b>int</b> x1,x2,y; x1 = (y*3)+x2-1; </pre> | <pre> t1 = 3 x1 = y * t1 t2 = x1 + x2 t3 = 1 x1 = t2 - t3 </pre> |
| El programador definió x1,x2,y                    | El compilador crea sus propias variables temporales (t1,t2,...)  |

### Las etiquetas son globales

Es decir, las etiquetas son únicas.

```

Proc Lectura
  E1: x = 10
      ⋮
Proc Calculo
  E1: t1 = 5

```

**Error!**, si ambas pertenecen al mismo programa, se debe poner distinta enumeración

Ejemplo: El triangulo, Leer N (Validar que sea  $>0$ ) y mostrar en consola un triangulo de N lineas de asteriscos, e.g. N = 4

```

*
* *
* * *
* * * *

```

### *Solución*

|   |  |
|---|--|
| Proc Lectura                                | Proc Linea                             |
| t1 = 0                                      | i=1                                    |
| E1: Writes ("Cantidad_de_Lineas? $\beta$ ") | E2: t1=(i $\leq$ k)                    |
| Read(N)                                     | <b>if</b> (t2=0) $\Rightarrow$ Goto E3 |
| t2 = (N $\leq$ t1)                          | Writes ("*")                           |
| <b>if</b> (t2=1) $\Rightarrow$ Goto E1      | INC i                                  |
| RET   | Goto E2                                |
|   | E3: NL                                 |
|   | RET                                    |

```

$MAIN
    CALL Lectura
    k = 1
E4:  t1=(k ≤ N)
      if(t1=0) ⇒ Goto E5
    CALL Linea
    INC k
    Goto E4
E5:  RET

```

## 2.4. Esquema de Traducción

Consiste en estrategias que el compilador debe seguir para realizar la traducción de las construcciones de programación (if, while, for, ...)

### 2.4.1. Las expresiones (booleanas y aritméticas)

```

/* expr = expresiones aritméticas
exprBoole = expresiones booleanas */

```

Las expresiones tradicionales por el compilador usando un algoritmo ideado por cada diseñador. En muchos casos, se aplica de la optimización.

Por este motivo, dejamos la traducción de las expresiones al libre albedrío.

Si anotamos  $\rightarrow \text{C3-Expr}(t_i)$   
 //  $t_i$  = temporal i-esima

Es el código-3 de la expr, cuyo resultado está en el temporal  $t_i$

```

/* Usamos  $t_i$  y no  $t_2$ ,  $t_3$  o  $t_4$ ... porque no sabemos cual es el número de la temporal final */

```

(o) Por ejemplo

Se tiene la Expr1

$$y * 3 + 4$$

convertir a C3

*Solución*

```

t1 = 3
t2 = y * t1
t3 = 4
[ t4 ] = t2 + t3

```

$\rightarrow \text{C3-Expr1}(t_4)$  //El resultado está en  $t_4$

(o) Convertir a C3 la siguiente `exprBoole1`

$$(x \leq 0) \text{ and } z$$

*Solución*

```
t1 = 0
t2 = x ≤ t1      → C3 - exprBoole1(t3) //El resultado esta en t3
[ t3 ] = t2 and z
```

### 2.4.2. La Asignación

En un lenguaje de Alto nivel

$$x = \text{expr};$$

Escribir un esquema de traducción C3 para esta construcción.

*Solución*

```
C3 - Expr(ti)
x = ti
```

Por ejemplo, usando el esquema de traducción anterior, Convertir a C3.

$$z = 2 * y - x$$

*Solución*

```
t1 = 2
t2 = t1 * y      → C3 - Expr(t3)
[ t3 ] = t2 - x
z = t3           → var = t3
```

### 2.4.3. Traducción para el If y While

| Java                          |               | C-3  |
|-------------------------------|---------------|--|
| <code>if (exprBoole) {</code> |               | <code>C3-ExprBoole( ti )</code>                            |
| <code>Sentencia ;</code>      | $\rightarrow$ | <code>if( ti = 0 ) <math>\Rightarrow</math> Goto Ef</code> |
| <code>}</code>                |               | <code>C3 - Sentencia</code>                                |
|                               |               | <code>Ef:</code>   |

(o) Ejemplos:

| Java  |   | C-3   |
|---|---|---|
| <b>if</b> (x<z){<br>x = z + x;<br>}                       | → | t1 = (x < z)<br><b>if</b> (t1 = 0) ⇒ Goto E1<br>t2 = z + x<br>x = t2<br>E1:   |
| <b>if</b> (x>0){<br><b>if</b> (z == 5){<br>z++;<br>}<br>} | → | t0 = 0<br>t1 = (x > t0)<br><b>if</b> (t1 = 0) ⇒ Goto E1<br>t2 = 5<br>t3 = (z == 5) ⇒ Goto E2<br>INC z<br>E2:<br>E1: |

(o) Escriba el esquema de traducción para el if-else

| Java  |   | C-3   |
|---|---|---|
| <b>if</b> (exprBoole){<br>Sentencia1;<br>}<br><b>else</b> {<br>Sentencia2;<br>} | → | C3-ExprBoole( ti )<br><b>if</b> ( ti = 0 ) ⇒ Goto Ef1<br>C3-Sentencia1<br>Goto Ef2<br>Ef1 :<br>C3-Sentencia2<br>E2: |

(o) Para el while

| Java   |   | C-3  |
|--|---|--|
| <b>while</b> (ExprBoole){<br>Sentencia;<br>} | → | Ef0 : C3-ExprBoole( ti )<br><b>if</b> ( ti=0 ) ⇒ Goto Ef<br>Sentencia<br>Goto Ef0<br>Ef: |

## 2.5. Representación del Código-3

Cada instrucción del C-3 se representa con un:

Block (Registros, Record, Struct, Class con todos los campos public)

Nota.- Los campos de este Block son de tipos primitivos (enteros)

La cantidad de campos que se utiliza bautiza la representación (de nombre)



- Terceto = 3 campos → Menos memoria, algoritmos mas complicados
- Cuadrupla = 4 campos → Cant. media de memoria, algoritmos termino medio (la asignada para la V.M.)
- Quintupla = 5 campos → Mucha memoria, algoritmos muy simples (Deprecated)

```
Class Cuadrupla{
    public int Op_code; //I.L. de Java
    public int Dir1, Dir2, Dir3; //Direcciones
}
```

Graficamente:

Cuadrupla = 

|         |      |      |      |
|---------|------|------|------|
| Op_code | Dir1 | Dir2 | Dir3 |
|---------|------|------|------|

//Op\_code = **Operation - Code**

Obviamente, no todas las instrucciones usan las tres direcciones. En estos casos, simplemente el campo **Dir** no utilizado se ignora (I don't care) y se denota con un underscore (“\_”)

Ejemplo:

$x = 15 \rightarrow$ 

|              |       |    |   |
|--------------|-------|----|---|
| OpAsignacion | Dir x | 15 | _ |
|--------------|-------|----|---|

El que diseña el compilador, crea el I.L., Una vez definidas las operaciones del C-3, usualmente define constantes para el “Opcode”.

| En Delphi   | En Java   | El programador hace |   |
|---|---|---------------------|---|
| Const<br>opMAS = 1;<br>opMENOS = 2;<br>opAND = 3; | Public <b>static final</b> opMAS = 1;<br>Public <b>static final</b> opMENOS = 2;<br>(etc) | OpAND               |   |
|   |   | ↓                   | El compilador lo sustituye por su valor |
|   |   | 3                   |   |

La inserción de los campos en la cuadrupla se leen de Izquierda a Derecha, tal como lo muestra textualmente la instrucción

Ejemplo:

|                        |          |   |   |   |           |
|------------------------|----------|---|---|---|-----------|
| $x = y \text{ and } z$ | OpAND    | x | y | z | Operación |
| $z = -x$               | OpMINUS  | z | x | _ | Operación |
| E1:                    | Etiqueta | 1 | _ | _ | Etiqueta  |
| Goto E2:               | OpGoto   | 2 | _ | _ | Jump      |

Al decir **Op** es que va a efectuar algo la maquina, no ponerlo a las etiquetas!

El **IF** que se utilizamos en C-3 no es un **If** de verdad, solo se usa para representación. En el salto se recomienda tener dos Opcodes (**IF var = x → Goto E**) **OpIF0** y **OpIF1**

Ejemplo:

|  |       |   |   |   |
|--|-------|---|---|---|
| <b>IF</b> <b>z</b> = 0 ⇒ <b>Goto</b> <b>E5</b> | OpIF0 | z | 5 | — |
| <b>IF</b> <b>z</b> = 1 ⇒ <b>Goto</b> <b>E9</b> | OpIF1 | z | 9 | — |

5 y 9 Son las Etiquetas **E5** y **E9** respectivamente.

Entonces un programa en C-3 es una Collection (Vector, Lista, Archivo,...) de **Cuadru-  
plas**.

Por Ejemplo:

|   |   |             |                |          |          |
|---|---|-------------|----------------|----------|----------|
| <b>x</b> = 27                                     | 0 | OpAssignNum | <b>x</b>       | 27       | —        |
| <b>y</b> = 1                                      | 1 | OpAssignNum | <b>y</b>       | 1        | —        |
| <b>p</b> = <b>x</b> and <b>y</b>                  | 2 | OpAnd       | <b>p</b>       | <b>x</b> | <b>y</b> |
| <b>if</b> ( <b>p</b> = 0) ⇒ <b>Goto</b> <b>E1</b> | 3 | OpIF0       | <b>p</b>       | 1        | —        |
| <b>writes</b> ("Hola")                            | 4 | OpWrites    | Dir. de "Hola" | —        | —        |
| <b>E1</b> :                                       | 5 | Etiqueta    | 1              | —        | —        |

## 2.6. Tabla de Simbolos

/\* Tabla de Simbolos = T.S. = Symbol Table = S.T. \*/

Tradicionalmente, se dice tabla de símbolos a un conjunto de tablas, no a una sola. Pero, para evitar esa ambigüedad, actualmente se llama: “**Administrador de Tablas de Símbolos**”. Para V.M. (Virtual Machine) que se quiere implementar, utilizaremos dos Tablas de Símbolos.

**TSID** = Tabla de Identificadores

**TSS** = Tablas de String's ctte /\*Println("Hola"); Hola → String ctte\*/

### 2.6.1. La T.S.I.D.

La **TSID** utiliza tuplas cuyos atributos son:

|          |        |        |         |
|----------|--------|--------|---------|
| NombreID | ValorI | ValorD | CantTMP |
|----------|--------|--------|---------|

CantTMP = Cantidad de Temporales

En esta tabla se almacenan las variables y los procedimientos.

- 1) Si la tupla es una variable:
  - ValorF contiene el tipo (es siempre negativo, Integer = -2, Boolean = -3)
  - ValorI mantiene el valor actual de la variable
  - CantTMP = No se usa
- 2) Si la tupla es un Procedimiento:
 

ValorI y ValorF almacenan, respectivamente el inicio y el fin de su C3
- 3) Para el campo ValorF se puede saber si la tupla es un Procedimiento o Variable:
  - $\text{ValorF} < -1 \Rightarrow \text{Variable}$
  - $\text{ValorF} \geq -1 \Rightarrow \text{Procedimiento}$

*Ejemplo:*

| TSID |           |        |        |         |
|------|-----------|--------|--------|---------|
|      | Nombre    | ValorI | ValorF | CantTMP |
| 0    | "x"       | ?      | -2     | —       |
| 1    | "Lectura" | 0      | 2      | 2       |
| 2    | "y"       | ?      | -3     | —       |
| 3    | "\$Main"  | 3      | 6      | 1       |

int x;  
boolean y;

| Codigo-3 (Vector de Cuadruplas) |      |         |  |  |
|---------------------------------|------|---------|--|--|
| 0                               | —    |         |  |  |
| 1                               | —    |         |  |  |
| 2                               | Ret  |         |  |  |
| 3                               | Call | Lectura |  |  |
| 4                               | —    |         |  |  |
| 5                               | —    |         |  |  |
| 6                               | Ret  |         |  |  |

Cuadrupla C3 [ ]

### 2.6.2. La T.S.S.

// T.S.S. = Tabla de String Ctte

Es simplemente una tabla de un solo atributo y almacena los string's cttes del programa

| T.S.S. |                    |
|--------|--------------------|
| 0      | "Hola Mundo"       |
| 1      | "*"                |
| 2      | "El factorial es " |
| 3      | "bye"              |

### 2.6.3. Codificaciones de las direcciones

Cuadrupla = 

|        |      |      |      |
|--------|------|------|------|
| OpCode | Dir1 | Dir2 | Dir3 |
|--------|------|------|------|

Si se refiere a un temporal ( $t_1, t_2, \dots$ ) o una etiqueta, simplemente se le anota su numero.

$$\begin{aligned}
 t2 = t4 + t7 &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpMas} & 2 & 4 & 7 \\ \hline \end{array} \\
 &\quad \quad \quad // \quad t2 \quad t4 \quad t7 \\
 \text{Goto E9} &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpGoto} & 9 & \_ & \_ \\ \hline \end{array} \\
 &\quad \quad \quad // \quad \text{E9} \\
 \text{if (t4 = 0)} \Rightarrow \text{Goto E5} &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpIf0} & 4 & 5 & \_ \\ \hline \end{array} \\
 &\quad \quad \quad // \quad t4 \quad \text{E5}
 \end{aligned}$$

Si se refiere a un identificador (ID) (Variables o procedimientos definidos por el programador), se usa el indice de la tabla negativo

| T.S.I.D. |        |  |  |  |
|----------|--------|--|--|--|
|          | Nombre |  |  |  |
| 0        | "x"    |  |  |  |
| 1        | "Base" |  |  |  |
| 2        | "Fact" |  |  |  |

$$\begin{aligned}
 x = t4 + \text{Base} &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpMas} & 0 & 4 & -1 \\ \hline \end{array} \\
 \text{Call Fact} &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpCall} & -2 & \_ & \_ \\ \hline \end{array} \\
 t2 = \text{Base} - \text{Base} &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpMenos} & 2 & -1 & -1 \\ \hline \end{array}
 \end{aligned}$$

Porque 0, -1, -2 ? son los indices en el **T.S.I.D** de dichas variables/procedimientos, si es positivo es un temporal o una etiqueta. **No olvides que las OpCodes son una función, no una Cadena!**

De la misma forma se hace con las String Ctte. (Operación: write o writes). Por ejemplo, Usando la **T.S.S.** del principio:

$$\begin{aligned}
 \text{writes}(\text{"Bye"}) &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpWrites} & -3 & \_ & \_ \\ \hline \end{array} \\
 \text{write}(\text{Base}) &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpWrite} & -1 & \_ & \_ \\ \hline \end{array} \\
 \text{writes}(\text{"*"}) &\Rightarrow \begin{array}{|c|c|c|c|} \hline \text{OpWrites} & -1 & \_ & \_ \\ \hline \end{array}
 \end{aligned}$$

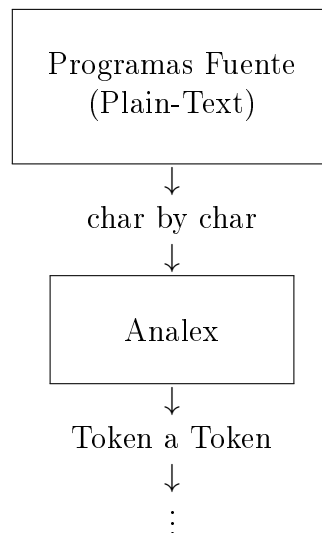
# Capítulo 3

## Analizador Léxico

*/\* Tambien se lo conoce como scanner, lexer o Analex \*/*

El analizador es el encargado de:

- “Descompone” el programa fuente en piezas léxicas, llamado Tokens
- “Elimina” los comentarios del programa fuente
- “Instala” (insert) los string cttes. a la TSS
- Si el lenguaje es no tipado (no tiene tipos) instala los ID’s a la TSID. (Si el lenguaje es tipado, la instalación de ID’s lo hace Parser)



*/\* Plain-Text = Texto (archivo) plano = Que se puede leer en el block de Notas \*/*

**Los Caracteres Char:** En un principio cada computadora definía los chars a su antojo:

| Computadora 1 | Computadora 2 |
|---------------|---------------|
| A = 0         | A = 2         |
| B = 1         | B = 3         |
| C = 2         | C = 4         |
| D = 3         | D = 5         |

La Computadora 1 envía el text “DDC” a la Computadora 2

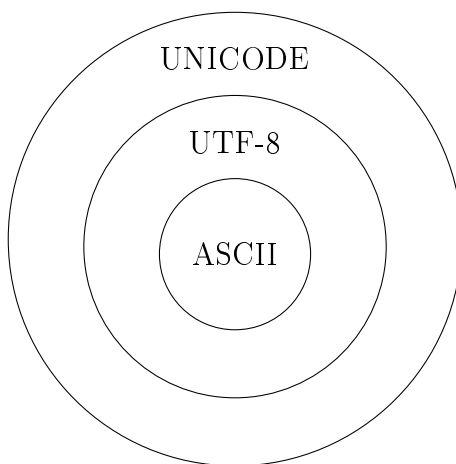
$$\text{Comp. 1} \xrightarrow{\text{DDC (3,3,2)}} \text{Comp. 2} \rightsquigarrow \text{recibe “BBA”}$$

Para evitar esto, crearon:

- Un código estándar de 7 bits llamado A.S.C.I.I. (Codigo de 0 a 127)

0 - 31 = Char no imprimibles  
 62 = “Space”  
 ⋮  
 65 = A  
 66 = B’  
 ⋮

- UTF-8 los char van desde 0 al 255
- Unicode (0 - 66535)

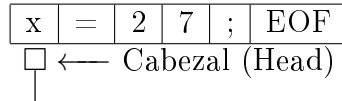


ASCII = 7 bits, UTF = 1 byte, UNICODE = 2 bytes

### 3.1. Manejo del Programa Fuente

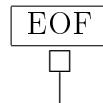
Para la manipular el programa se utiliza un ADT (Maquina) llamado “**Cinta de Caracteres**” o multilíneas o serializador

Para el Analex, la cinta se muestra así:



Cada celda de la cinta aloja un char (carácter) y siempre la cinta termina en un char especial EOF (End Of File)

Por ejemplo una cinta vacia:

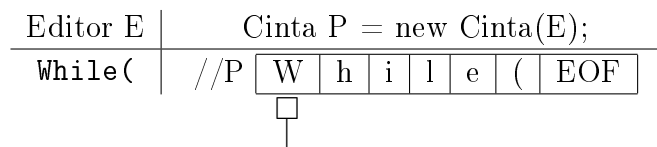


El Cabezal o Head permite leer el caracter o char que esta “mirando”. La cinta puede trabajar con un editor, un archivo, un puerto de conexión en la red, etc.

### Métodos de la Class Cinta

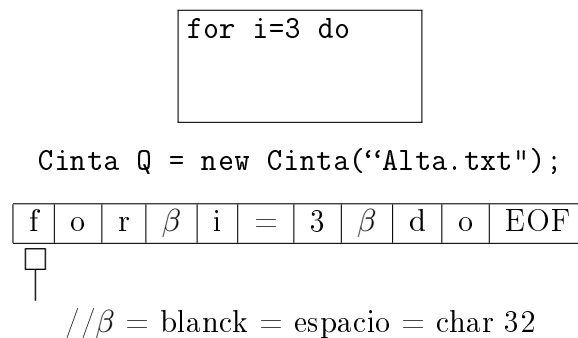
- `Public Cinta (Editor E)`
- `Public Cinta (String FileName)`
- `Public Cinta (Net algo)`

*Por ejemplo:*



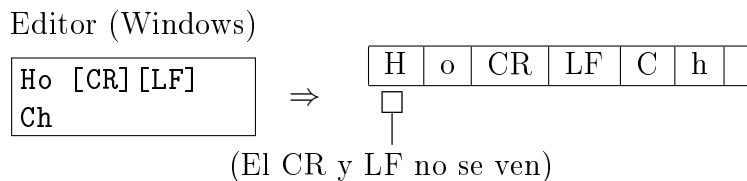
- `Public void init() //Rewind`  
   //Mueve el cabezal a la primera celda
- `Public void Avanzar() //Forward`  
   // Mueve el cabezal a la sgte. celda de la derecha
- `Public char cc() //Current-Char (Caracter del cabezal)`  
   //Devuelve el char que esta viendo el cabezal

*Por ejemplo:* Usando la Cinta: `Alfa.txt`



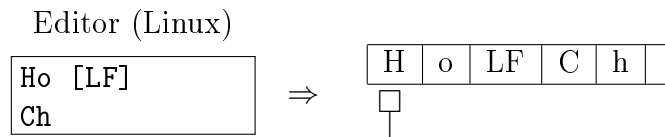






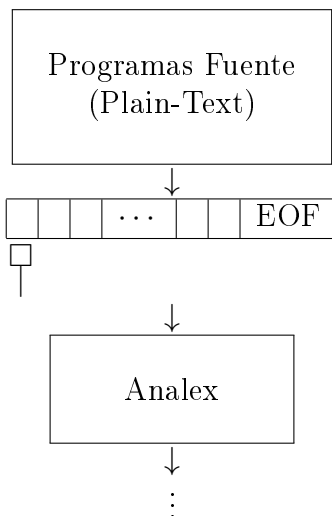
### Representación en Linux

Solo utiliza el LF



Para la cinta no tiene problema en trabajar en ambos formatos. Sin importar el programa fuente, fuera hecho en Linux o Windows, la cinta representa al salto de linea con un solo char (de caracteres): EOLN

Con la cinta, la Arquitectura del compilador, empieza así:



## 3.2. Especificación de los Componentes Lexicos

//Componentes Léxicos = Tokens

Matematicamente, un token es un lenguaje regular que no contiene a la cadena vacía ( $\lambda$ )

Un alfabeto define los char's (letras o símbolos) que se van a utilizar. Un lenguaje es un conjunto de cadenas que se forma con los char's de  $\Sigma$

Por Ejemplo:

$$\Sigma = \{a, b, c\} \text{ y } L = \{\lambda, "ab", "ccb"\}$$

Un lenguaje es regular si puede ser especificado por una expresión regular (REGEX).

### 3.2.1. Definiciones Regulares

$$Sol = a(a|e|i|o|u)^*b \text{ ó}$$

$$\begin{aligned} algo &= (a|e|i|o|u) \\ sol &= a(algo)^*b \end{aligned}$$

Permite especificar REGEX por partes, a traves de producciones. Una definición Regular no necesita especificar el alfabeto  $\Sigma$

#### Operadores

- $|$  = lease “o” (Unión)
- $.$  = concatenación
- $*$  = 0 o mas (Estrellas de Kleene)
- $+$  = 1 o mas (Estrella de Kleene)
- $\rightarrow$  = lease “es” o “puede ser”

#### Ejemplos:

(o) Especificar el Token NUM (Números enteros sin signo)

$$NUM = \{ "0", ..., "9", "3564", "800", etc \}$$

*Solución*

$$\text{Digitos} \rightarrow 0|1|2|3|4|5|6|7|8|9$$

$$NUM \rightarrow \text{Digitos}^+$$

(o) Especificar el Token NumP (para los números enteros pares)

$$NUM = \{ "0", ..., "22", "154", "430", etc \}$$

*Solución*

$$\text{Par} \rightarrow 0|2|4|6|8$$

$$\text{Digitos} \rightarrow 0|1|2|3|4|5|6|7|8|9$$

$$NUM \rightarrow \text{Digitos}^* \text{Par}$$

**Recuerda:** Los números (Enteros o Float) en compilación se anotan sin signo, -34 o +34 saldría ERROR

(o) Especificar el Token ID, el cual esta compuesto por solamente letras

$$ID = \{x, y, Area, Base, Altura, \dots\}$$

Letra  $\rightarrow a|b|c|\dots|z$

ID  $\rightarrow \text{Letra}^+$

/\* Cuando se anotan las letras en una definición regular (DEFREG), no es necesario escribir mayúsculas y minúsculas. Solo basta 1 de estas cosas.

UpperCase = Mayúsculas

LowerCase = Minúsculas

El compilador, según la política del lenguaje, hará una diferencia en los Case

- Si el lenguaje es Case-Sensitive, distinguirá las mayúsculas de las minúsculas. Por ejemplo: Java, C, C++, C# (`int a,A`; se refiere a distintos casos)
- Si el lenguaje no es Case-Sensitive, el compilador lleva todo a mayúsculas y trabaja internamente así. Por ejemplo: Pascal (`Var a,A: Integer`; se refiere al mismo ID)

\*/

(o) Escribir una DEFREG para el Token IDX, el cual esta formado con solo letras, pero tienen dos formas:

Empieza con: A *algo* Z

Empieza con: Z *algo* A

*Solución*

Letra  $\rightarrow a|b|c|\dots|z$

IDX  $\rightarrow a \text{ Letra}^* z \mid z \text{ Letra}^* a$

### 3.2.2. Lexemas (Lexem's)

Simplemente, se llama lexemas a los Strings que pertenecen a un Token:

$$\begin{array}{ccccccc} \text{NUM} & = & \{ \dots & 315, & 860, & 1001, & \dots \} \\ & & \uparrow & \nwarrow & \uparrow & \nearrow & \\ & & \text{Token} & & \text{Lexemas} & & \end{array}$$

Del Ejercicio de DEFREG, ZBAA es un lexema del token IDX

### 3.2.3. Números Reales

En compilación un número entero no es un número real (flotante)

$$2,0 \neq 2$$

Todo numero real si o si debe tener el **punto decimal** (.), si no lo tiene se considera entero.

### Números Reales Formales

Los números reales formales son aquellos que usan las siguientes 3 partes:

Parte entera Punto(.) Parte decimal

Ejemplo: 3.17 (V)     .24 (X)  $\rightarrow$  0.24(V)

### Números Reales Informales

Se puede obviar la `Parte entera` o la `Parte decimal` (pero no ambas)

Ejemplo: 35. (V)     .24 (V)

(o) Mediante una Def. regular, especificar. Los NUMR Formales (NUMR = Números reales):

*Solución*

`Digito`  $\rightarrow$  0|1|2|...|9

`NUMR`  $\rightarrow$  `Digito`<sup>+</sup> . `Digito`<sup>+</sup>

**Nota:** Dado que las producciones `Digitos`  $\rightarrow$  y `Letra`  $\rightarrow$ , son muy utilizadas, no es necesario definir las todo el tiempo

## 3.3. Diagrama de Transiciones (dt)

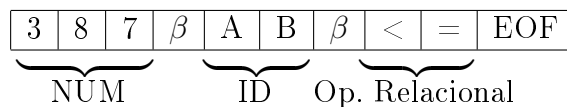
Puesto que los Tokens son Lenguajes Regulares, por un teorema, es posible reconocer utilizando un AFD (Autómata Finito Determinístico). Pero... Un compilador asume que el alfabeto ( $\Sigma$ ) es infinito.

Para solucionar esto usamos un AFD que use definiciones regulares, llamado **diagrama de transiciones o dt**.

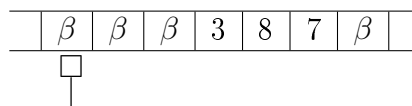
Las definiciones regulares permitidas en un dt solo agrupan **char's**. Es decir, no usan la Estrella de Kleene ni la concatenación.

### 3.3.1. El ANALEX como un jugador de sopa de letras

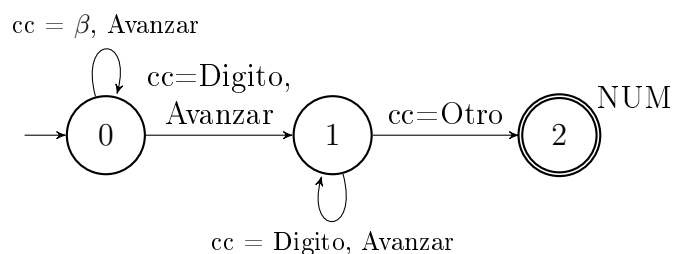
El ANALEX busca en la cinta (SCAN) los lexemas que pertenecen a un Token, por ejemplo:



(o) En una cinta se encuentra un NUM dibuja un dt para reconocerlo. (NUM = números enteros sin signo)



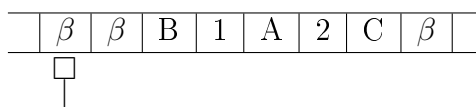
*Respuesta:*



//cc = Caracter del Cabezal

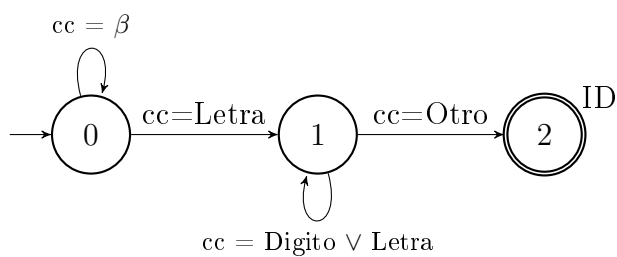
Cuando pasas de (1) a (2) gracias a “Otro” se dice: **lo que esta detrás de mi es un Número.**

(o) En una cinta hay un ID que empieza con letra y luego le sigue una combinación de letras y dígitos.



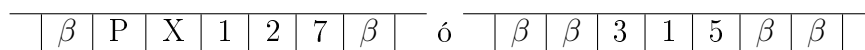
Dibujar un dt para reconocer este Token.

*Respuesta:*

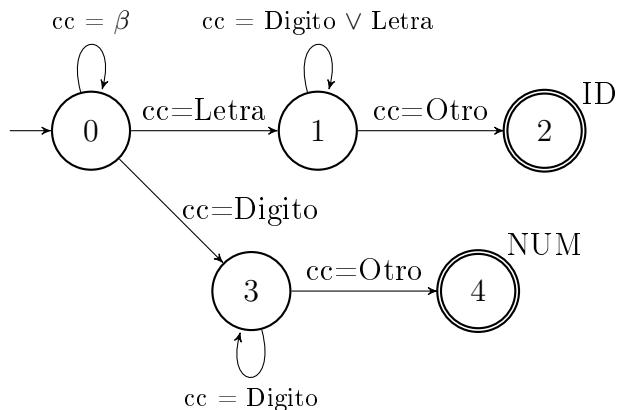


Cuando pasas de (1) a (2) gracias a “Otro” se dice: **lo que esta detrás de mi es un ID.**

(o) En una cinta hay un NUM o un ID, hacer un dt para reconocerlo.



*Solución:*



**Simplificación:** Es posible eliminar “cc=” de las aristas

### 3.3.2. El Espacio

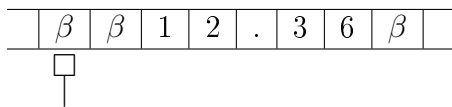
El espacio se usa para separar las palabras en el programa fuente. Según la Norma ASCII. El compilador toma como espacio a **tres caracteres**.

- $\beta$  = barra espaciadora = ASCII 32
- TAB = tabulador = ASCII 9
- EOLN = fin de línea = LF = ASCII 10

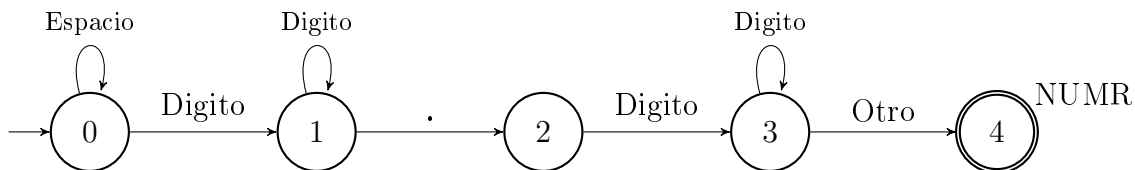
Entonces en vez de usar “ $\beta$ ” usamos la producción

**Espacio**  $\rightarrow \beta | \text{TAB} | \text{EOLN}$

(o) Dibujar un dt que reconozca NUMR formal que esta en la Cinta



*Solución:*



### 3.3.3. Los Tokens ERROR y FIN

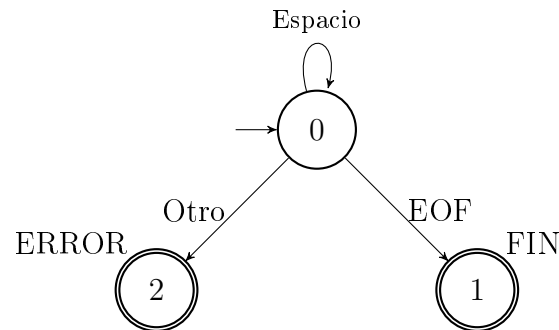
#### El Token ERROR

El dt, devuelve el Token ERROR cuando no es posible reconocer un Token

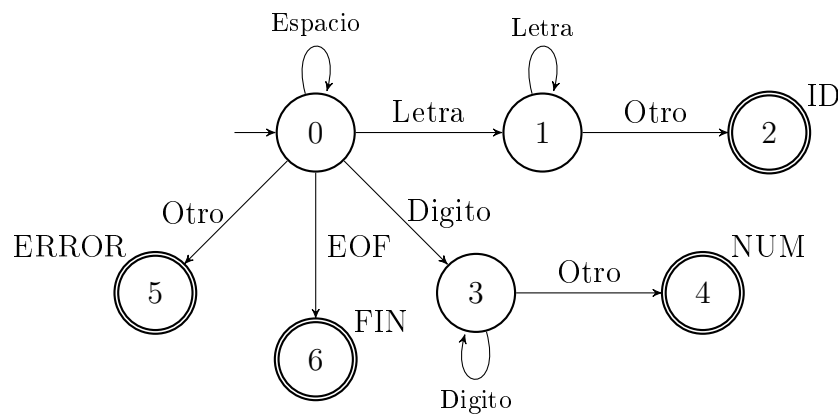
### El Token Fin

Este Token es devuelto por el dt, cuando el **cc** alcanza al EOF.

La estructura general seria:



(o) Dibujar un dt que reconozca NUM (enteros sin signo) e ID (solo letras)



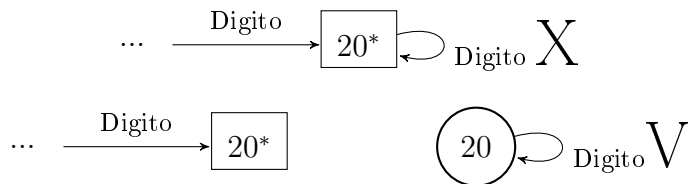
#### Recuerda:

- Todas las transiciones implican un **Avanzar**, excepto “Otro” y “EOF”
- **Espacio** =  $\beta$ |TAB|EOLN
- Uno puedo llamar al dt cuantas veces quiera
- Al llegar a un estado aceptado se dice: “Lo que esta detrás de mi es un ID, NUM, NUMR...”

### 3.4. “Reglas” de los dt

- 1) Existe 1 y solo 1 estado inicial (tradicionalmente el estado inicial es el 0)
- 2) Todos los estados son diferentes.

Error Común: Cuando el dt tiene muchos estados y aristas, el diseñador utiliza conectores (los conectores no son estados)



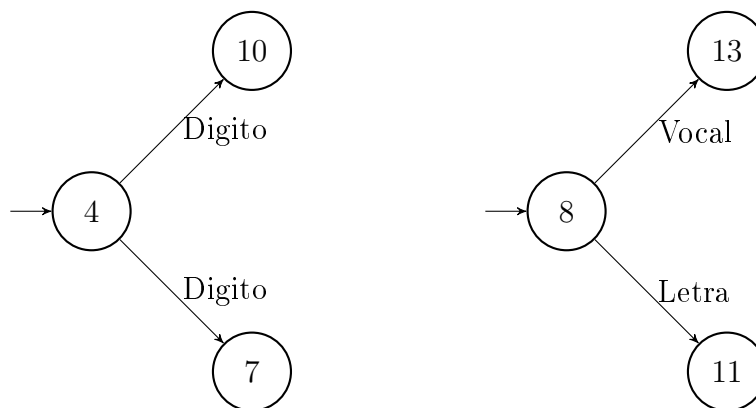
- 3) Si un estado tiene aristas salientes entonces tiene una arista rotulada con “Otro”.

Errores Comunes:

- Usar “Otro” como una única arista
- Tener dos “Otro” saliendo de un estado
- Usar “Otro” como ciclo de un estado

- 4) Un dt no puede presentar un No-Determinismo (ambigüedad en las transiciones)

Errores Comunes:



No puede haber dos aristas que compartan el mismo conjunto (o subconjunto) de caracteres ( $\text{Digito} = \text{Digito}$ ,  $\text{Vocal} \subseteq \text{Letra}$ )

Para solucionarlo se debe especificar en las definiciones regulares las diferencias que tienen entre si:

**Vocal** =  $a|e|i|o|u$

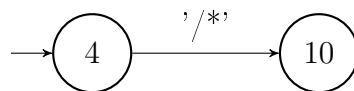
**Letra** =  $b|c|d|f|\dots|z$  // Letras sin vocales

- 5) El dt debe tener al menos un estado final. Si no tiene un estado final, nunca finalizara.

Errores Comunes:

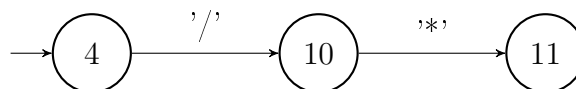
- Definir mal una definición regular. **Recuerda:** el dt analiza char a char el cual le otorgan las definiciones regulares ( $a|b|c|\dots|z$ ) estas no deben agruparse en mas de dos ( $ab|cd|dd|\dots|za$ ) puesto que cc (Cabeza de la Cinta) no lo reconocería
- Poner mas de un char en la arista





Error: Hay dos char's

La solución sería crear un nuevo estado que separe ambos caracteres:



(o) Dibujar un dt que reconozca Tokens usando solamente letras

|  |   |   |         |   |   |   |         |   |  |
|--|---|---|---------|---|---|---|---------|---|--|
|  | C | A | $\beta$ | A | E | Z | $\beta$ | B |  |
|--|---|---|---------|---|---|---|---------|---|--|

IDV  $\rightarrow$  termina en Vocal (CA,...)

ID  $\rightarrow$  no termina en Vocal (AEZ, B,...)

*Solución:*

Const = B|C|D|...|Z //no tiene vocales

Vocal = A|E|I|O|U

