

Análisis Léxico y Sintáctico

Casiano R. León ¹

22 de mayo de 2012

¹DEIOC Universidad de La Laguna

Índice general

1. Expresiones Regulares en sed	6
1.1. Transferencia de Control	6
1.2. Inserción de Texto	7
1.3. Trasferencia de Control Condicional	8
1.4. Rangos	9
1.5. Uso de la negación	11
1.6. Siguiendo Línea: La orden <i>n</i>	12
1.7. Manipulando tablas numéricas	14
1.8. Traducción entre Tablas	15
1.9. Del espacio de Patrones al de Mantenimiento	15
1.10. La orden <i>N</i>	18
1.11. Suprimir: El Comando <i>D</i>	20
1.12. Búsqueda entre líneas	21
1.13. Seleccionando Items en un Registro Multilínea	22
 2. Expresiones Regulares en Flex	 24
2.1. Estructura de un programa LEX	24
2.2. Versión Utilizada	27
2.3. Espacios en blanco dentro de la expresión regular	28
2.4. Ejemplo Simple	28
2.5. Suprimir	29
2.6. Declaración de <code>yytext</code>	29
2.7. Declaración de <code>yylex()</code>	30
2.8. <code>yywrap()</code>	31
2.9. <code>unput()</code>	32
2.10. <code>input()</code>	33
2.11. <code>REJECT</code>	34
2.12. <code>yyomore()</code>	34
2.13. <code>yyless()</code>	34
2.14. Estados	35
2.15. La pila de estados	38
2.15.1. Ejemplo	38
2.16. Final de Fichero	39
2.17. Uso de Dos Analizadores	40
2.18. La Opción <code>outfile</code>	41
2.19. Leer desde una Cadena: <code>YY_INPUT</code>	41
2.20. El operador de “trailing context” o “lookahead” positivo	42
2.21. Manejo de directivas <code>include</code>	43
2.22. Análisis Léxico desde una Cadena: <code>yy_scan_string</code>	46
2.23. Análisis de la Línea de Comandos y 2 Analizadores	47
2.24. Declaraciones <code>pointer</code> y <code>array</code>	51
2.25. Las Macros <code>YY_USER_ACTION</code> , <code>yy_act</code> e <code>YY_NUM_RULES</code>	52
2.26. Las opciones <code>interactive</code>	53

2.27. La macro <code>YY_BREAK</code>	53
--	----

3. Expresiones Regulares en Perl 56

3.1. Introducción	56
3.1.1. Un ejemplo sencillo	57
3.1.2. Depuración de Expresiones Regulares	82
3.1.3. Tablas de Escapes, Metacaracteres, Cuantificadores, Clases	83
3.1.4. Variables especiales después de un emparejamiento	87
3.1.5. Ambito Automático	90
3.1.6. Opciones	91
3.2. Algunas Extensiones	93
3.2.1. Comentarios	93
3.2.2. Modificadores locales	93
3.2.3. Mirando hacia adetrás y hacia adelante	94
3.2.4. Definición de Nombres de Patrones	102
3.2.5. Patrones Recursivos	104
3.2.6. Cuantificadores Posesivos	114
3.2.7. Perl 5.10: Numeración de los Grupos en Alternativas	116
3.2.8. Ejecución de Código dentro de una Expresión Regular	117
3.2.9. Expresiones Regulares en tiempo de matching	122
3.2.10. Expresiones Condicionales	126
3.2.11. Verbos que controlan el retroceso	130
3.3. Expresiones Regulares en Otros Lenguajes	134
3.4. Casos de Estudio	138
3.4.1. Secuencias de números de tamaño fijo	138
3.4.2. Palabras Repetidas	140
3.4.3. Análisis de cadenas con datos separados por comas	141
3.4.4. Las Expresiones Regulares como Exploradores de un Árbol de Soluciones	143
3.4.5. Número de substituciones realizadas	147
3.4.6. Expandiendo y comprimiendo tabs	148
3.4.7. Modificación de Múltiples Ficheros: one liner	149
3.5. <code>tr</code> y <code>split</code>	149
3.6. <code>Pack</code> y <code>Unpack</code>	151
3.7. Práctica: Un lenguaje para Componer Invitaciones	152
3.8. Análisis Sintáctico con Expresiones Regulares Perl	153
3.8.1. Introducción al Análisis Sintáctico con Expresiones Regulares	153
3.8.2. Construyendo el AST con Expresiones Regulares 5.10	161
3.9. Práctica: Traducción de <code>invitation</code> a HTML	169
3.10. Análisis Sintáctico con <code>Regexp::Grammars</code>	171
3.10.1. Introducción	171
3.10.2. Objetos	179
3.10.3. Renombrando los resultados de una subregla	180
3.10.4. Listas	182
3.10.5. Pseudo sub-reglas	190
3.10.6. Llamadas a subreglas desmemoriadas	193
3.10.7. Destilación del resultado	196
3.10.8. Llamadas privadas a subreglas y subreglas privadas	200
3.10.9. Mas sobre listas	200
3.10.10. La directiva <code>require</code>	210
3.10.11. Casando con las claves de un hash	212
3.10.12. Depuración	214
3.10.13. Mensajes de <code>log</code> del usuario	219
3.10.14. Depuración de Regexp	220

3.10.15.	Manejo y recuperación de errores	221
3.10.16.	Mensajes de Warning	224
3.10.17.	Simplificando el AST	224
3.10.18.	Reciclando una <code>Regex::Grammar</code>	228
3.10.19.	Práctica: Calculadora con <code>Regex::Grammars</code>	237
4.	La Estructura de los Compiladores: Una Introducción	239
4.1.	Las Bases	239
4.1.1.	Repaso: Las Bases	243
4.1.2.	Práctica: Crear y documentar el Módulo <code>PL::Tutu</code>	244
4.2.	Las Fases de un Compilador	244
4.2.1.	Repaso: Fases de un Compilador	250
4.2.2.	Práctica: Fases de un Compilador	252
4.3.	Análisis Léxico	254
4.3.1.	Ejercicio: La opción <code>g</code>	256
4.3.2.	Ejercicio: Opciones <code>g</code> y <code>c</code> en Expresiones Regulares	257
4.3.3.	Ejercicio: El orden de las expresiones regulares	257
4.3.4.	Ejercicio: <code>Regex</code> para cadenas	257
4.3.5.	Ejercicio: El <code>or</code> es vago	258
4.3.6.	Práctica: Números de Línea, Errores, Cadenas y Comentarios	258
4.4.	Pruebas para el Analizador Léxico	262
4.4.1.	Comprobando el Analizador Léxico	265
4.4.2.	Práctica: Pruebas en el Análisis Léxico	269
4.4.3.	Repaso: Pruebas en el Análisis Léxico	277
4.5.	Conceptos Básicos para el Análisis Sintáctico	279
4.5.1.	Ejercicio	279
4.6.	Análisis Sintáctico Predictivo Recursivo	280
4.6.1.	Introducción	280
4.6.2.	Ejercicio: Recorrido del árbol en un ADPR	282
4.6.3.	Ejercicio: Factores Comunes	282
4.6.4.	Derivaciones a vacío	283
4.6.5.	Construcción de los conjuntos de Primeros y Siguietes	284
4.6.6.	Ejercicio: Construir los <i>FIRST</i>	284
4.6.7.	Ejercicio: Calcular los <i>FOLLOW</i>	285
4.6.8.	Práctica: Construcción de los <i>FIRST</i> y los <i>FOLLOW</i>	285
4.6.9.	Gramáticas <i>LL(1)</i>	286
4.6.10.	Ejercicio: Caracterización de una gramática <i>LL(1)</i>	286
4.6.11.	Ejercicio: Ambigüedad y <i>LL(1)</i>	287
4.6.12.	Práctica: Un analizador APDR	287
4.6.13.	Práctica: Generación Automática de Analizadores Predictivos	287
4.7.	Esquemas de Traducción	293
4.8.	Recursión por la Izquierda	294
4.8.1.	Eliminación de la Recursión por la Izquierda en la Gramática	294
4.8.2.	Eliminación de la Recursión por la Izquierda en un Esquema de Traducción	295
4.8.3.	Ejercicio	295
4.8.4.	Convirtiendo el Esquema en un Analizador Predictivo	295
4.8.5.	Ejercicio	296
4.8.6.	Práctica: Eliminación de la Recursividad por la Izquierda	296
4.9.	Árbol de Análisis Abstracto	297
4.9.1.	Lenguajes Árbol y Gramáticas Árbol	297
4.9.2.	Realización del AAA para Tutu en Perl	300
4.9.3.	AAA: Otros tipos de nodos	306
4.9.4.	Declaraciones	307

4.9.5. Práctica: Arbol de Análisis Abstracto	307
4.10. Análisis Semántico	308
4.10.1. Práctica: Declaraciones Automáticas	310
4.10.2. Práctica: Análisis Semántico	310
4.11. Optimización Independiente de la Máquina	311
4.11.1. Práctica: Plegado de las Constantes	314
4.12. Patrones Árbol y Transformaciones Árbol	314
4.12.1. Práctica: Casando y Transformando Árboles	319
4.13. Asignación de Direcciones	320
4.13.1. Práctica: Cálculo de las Direcciones	321
4.14. Generación de Código: Máquina Pila	323
4.15. Generación de Código: Máquina Basada en Registros	326
4.15.1. Práctica: Generación de Código	331
4.16. Optimización de Código	332
4.16.1. Práctica: Optimización Peephole	333
5. Construcción de Analizadores Léxicos	334
5.1. Encontrando los terminales mediante sustitución	334
5.2. Construcción usando la opción <code>g</code> y el ancla <code>G</code>	337
5.3. La clase <code>Parse::Lex</code>	337
5.3.1. Condiciones de arranque	341
5.4. La Clase <code>Parse::CLex</code>	342
5.5. El Módulo <code>Parse::Flex</code>	344
5.6. Usando <code>Text::Balanced</code>	345
6. RecDescent	348
6.1. Introducción	348
6.2. Orden de Recorrido del Árbol de Análisis Sintáctico	350
6.3. La ambigüedad de las sentencias <code>if-then-else</code>	354
6.4. La directiva <code>commit</code>	360
6.5. Las Directivas <code>skip</code> y <code>leftop</code>	361
6.6. Las directivas <code>rulevar</code> y <code>reject</code>	363
6.7. Utilizando <code>score</code>	364
6.8. Usando <code>autoscore</code>	365
6.9. El Hash <code>%item</code>	366
6.10. Usando la directiva <code>autotree</code>	367
6.11. Práctica	369
6.12. Construyendo un compilador para Parrot	369
6.13. Práctica	379
6.14. Práctica	379
7. Análisis LR	380
7.1. <code>Parse::Yapp</code> : Ejemplo de Uso	380
7.2. Conceptos Básicos	386
7.3. Construcción de las Tablas para el Análisis SLR	388
7.3.1. Los conjuntos de Primeros y Sigüientes	388
7.3.2. Construcción de las Tablas	389
7.4. El módulo Generado por <code>yapp</code>	393
7.5. Algoritmo de Análisis LR	396
7.6. Depuración en <code>yapp</code>	397
7.7. Precedencia y Asociatividad	398
7.8. Generación interactiva de analizadores <code>Yapp</code>	402
7.9. Construcción del Árbol Sintáctico	403
7.10. Acciones en Medio de una Regla	405

7.11. Esquemas de Traducción	406
7.12. Definición Dirigida por la Sintaxis	406
7.13. Manejo en yapp de Atributos Heredados	408
7.14. Acciones en Medio de una Regla y Atributos Heredados	412
7.15. Recuperación de Errores	414
7.16. Recuperación de Errores en Listas	416
7.17. Consejos a seguir al escribir un programa yapp	417
7.18. Práctica: Un C simplificado	421
7.19. La Gramática de yapp / yacc	425
7.19.1. La Cabecera	425
7.19.2. La Cabecera: Diferencias entre yacc y yapp	426
7.19.3. El Cuerpo	427
7.19.4. La Cola: Diferencias entre yacc y yapp	428
7.19.5. El Análisis Léxico en yacc : flex	429
7.19.6. Práctica: Uso de Yacc y Lex	430
7.20. El Analizador Ascendente Parse::Yapp	430
7.21. La Estructura de Datos Generada por YappParse.yp	434
7.22. Práctica: El Análisis de las Acciones	436
7.23. Práctica: Autoacciones	436
7.24. Práctica: Nuevos Métodos	438
7.25. Práctica: Generación Automática de Árboles	439
7.26. Recuperación de Errores: Visión Detallada	439
7.27. Descripción Eyapp del Lenguaje SimpleC	441
7.28. Diseño de Analizadores con Parse::Eyapp	443
7.29. Práctica: Construcción del AST para el Lenguaje Simple C	446
7.30. El Generador de Analizadores byacc	446
8. Análisis Sintáctico con Parse::Eyapp	450
8.1. Conceptos Básicos para el Análisis Sintáctico	450
8.2. Parse::Eyapp : Un Generador de Analizadores Sintácticos	452
8.3. Depuración de Errores	461
8.4. Acciones y Acciones por Defecto	470
8.5. Traducción de Infijo a Postfijo	474
8.6. Práctica: Traductor de Términos a GraphViz	475
8.7. Práctica: Gramática Simple en Parse::Eyapp	477
8.8. El Método YName y la Directiva %name	478
8.9. Construyendo el Arbol de Análisis Sintactico Mediante Directivas	482
8.10. La Maniobra de bypass	488
8.11. Salvando la Información en los Terminales Sintácticos	490
8.12. Práctica: Análisis Sintáctico	492
8.13. Podando el Arbol	493
8.14. Nombres para los Atributos	497
8.15. Bypass Automático	502
8.16. La opción alias de %tree	508
8.17. Diseño de Analizadores con Parse::Eyapp	511
8.18. Práctica: Construcción del Arbol para el Lenguaje Simple	514
8.19. Práctica: Ampliación del Lenguaje Simple	514
8.20. Agrupamiento y Operadores de Listas	514
8.21. El método str en Mas Detalle	526
8.22. El Método descendant	532
8.23. Conceptos Básicos del Análisis LR	532
8.24. Construcción de las Tablas para el Análisis SLR	535
8.24.1. Los conjuntos de Primeros y Siguietes	535

8.24.2. Construcción de las Tablas	536
8.25. Algoritmo de Análisis LR	540
8.26. Precedencia y Asociatividad	543
8.27. Acciones en Medio de una Regla	547
8.28. Manejo en <code>eyapp</code> de Atributos Heredados	548
8.29. Acciones en Medio de una Regla y Atributos Heredados	551
9. Análisis Semántico con <code>Parse::Eyapp</code>	554
9.1. Esquemas de Traducción: Conceptos	554
9.2. Esquemas de Traducción con <code>Parse::Eyapp</code>	555
9.3. Definición Dirigida por la Sintaxis y Gramáticas Atribuidas	575
9.4. El módulo <code>Language::AttributeGrammar</code>	577
9.5. Usando <code>Language::AttributeGrammars</code> con <code>Parse::Eyapp</code>	579
10. Transformaciones Árbol con <code>Parse::Eyapp</code>	584
10.1. Árbol de Análisis Abstracto	584
10.2. Selección de Código y Gramáticas Árbol	587
10.3. Patrones Árbol y Transformaciones Árbol	589
10.4. El método <code>m</code>	591
10.5. Transformaciones Árbol con <code>treereg</code>	594
10.6. Transformaciones de Árboles con <code>Parse::Eyapp::Treeregexp</code>	596
10.7. La opción <code>SEVERITY</code>	609
11. Análisis Sintáctico con <code>yacc</code>	612
11.1. Introducción a <code>yacc</code>	612
11.2. Precedencia y Asociatividad	615
11.3. Uso de <code>union</code> y <code>type</code>	619
11.4. Acciones en medio de una regla	619
11.5. Recuperación de Errores	623
11.6. Recuperación de Errores en Listas	626
12. Análisis de Ámbito	631
12.1. Análisis de Ámbito: Conceptos	631
12.2. Descripción <code>Eyapp</code> del Lenguaje SimpleC	640
12.3. Práctica: Construcción del AST para el Lenguaje Simple C	652
12.4. Práctica: Análisis de Ámbito del Lenguaje Simple C	652
12.5. La Dificultad de Elaboración de las Pruebas	662
12.6. Análisis de Ámbito con <code>Parse::Eyapp::Scope</code>	670
12.7. Resultado del Análisis de Ámbito	686
12.8. Usando el Método <code>str</code> para Analizar el Árbol	691
12.9. Práctica: Establecimiento de la relación uso-declaración	700
12.10 Práctica: Establecimiento de la Relación Uso-Declaración Usando Expresiones Regulares Árbol	701
12.11 Práctica: Estructuras y Análisis de Ámbito	701
13. Análisis de Tipos	705
13.1. Análisis de Tipos: Conceptos Básicos	705
13.2. Conversión de Tipos	709
13.3. Expresiones de Tipo en Simple C	710
13.4. Construcción de las Declaraciones de Tipo con <code>hnew</code>	710
13.5. Inicio de los Tipos Básicos	713
13.6. Comprobación Ascendente de los Tipos	714
13.7. Análisis de Tipos: Mensajes de Error	716
13.8. Comprobación de Tipos: Las Expresiones	719
13.9. Comprobación de Tipos: Indexados	721

13.10Comprobación de Tipos: Sentencias de Control	724
13.11Comprobación de Tipos: Sentencias de Asignación	724
13.12Comprobación de Tipos: Llamadas a Funciones	726
13.13Comprobación de Tipos: Sentencia RETURN	730
13.14Comprobación de Tipos: Sentencias sin tipo	733
13.15Ejemplo de Árbol Decorado	734
13.16Práctica: Análisis de Tipos en Simple C	735
13.17Práctica: Análisis de Tipos en Simple C	735
13.18Práctica: Análisis de Tipos en Simple C con Gramáticas Atribuidas	735
13.19Práctica: Sobrecarga de Funciones en Simple C	737
13.20Análisis de Tipos de Funciones Polimorfas	737
13.20.1.Un Lenguaje con Funciones Polimorfas	737
13.20.2.La Comprobación de Tipos de las Funciones Polimorfas	743
13.20.3.El Compilador	751
13.20.4.Un Algoritmo de Unificación	752
13.21Práctica: Inferencia de Tipos	755
14.Instrucciones Para la Carga de Módulos en la ETSII	756
15.Usando Subversion	757

Índice de figuras

4.1. El resultado de usar <code>perldoc Tutu</code>	243
7.1. NFA que reconoce los prefijos viables	388
7.2. DFA equivalente al NFA de la figura 7.1	390
7.3. Esquema de herencia de <code>Parse::Yapp</code> . Las flechas continuas indican herencia, las punteadas uso. La cla	
8.1. NFA que reconoce los prefijos viables	534
8.2. DFA equivalente al NFA de la figura 8.23.1	537

Índice de cuadros

7.1. Tablas generadas por yapp . El estado 3 resulta de transitar con \$	393
7.2. Recuperación de errores en listas	417
8.1. El nodo PLUS recoge el número de línea del terminal+	493
8.2. Tablas generadas por eyapp . El estado 3 resulta de transitar con \$	540
12.1. Representación de AST con str	693

A Juana

*For it is in teaching that we learn
And it is in understanding that we are understood*

Agradecimientos/Acknowledgments

I'd like to thank Tom Christiansen, Damian Conway, Simon Cozens, Francois Desarmenien, Richard Foley, Jeffrey E.F. Friedl, Joseph N. Hall, Tim Jennes, Andy Lester, Allison Randall, Randal L. Schwartz, Michael Schwern, Peter Scott, Sriram Srinivasan, Lincoln Stein, Dan Sugalski, Leopold Tötsch, Nathan Torkington and Larry Wall for their books and/or their modules. To the Perl Community.

Special thanks to Larry Wall for giving us Perl.

A mis alumnos de Procesadores de Lenguajes en el primer curso de la Ingeniería Superior Informática en la Universidad de La Laguna

Capítulo 1

Expresiones Regulares en sed

El editor `sed` es un editor no interactivo que actúa ejecutando los comandos (similares a los de `ex`) que figuran en el guión `sed` (*script sed*) sobre los ficheros de entrada y escribiendo el resultado en la salida estandar. Normalmente se llama en una de estas dos formas:

```
sed [opciones] 'comando' fichero(s)
sed [opciones] -f guion fichero(s)
```

Si no se especifican ficheros de entrada, `sed` lee su entrada desde la entrada estandar.

Todos los comandos en el guión son ejecutados sobre todas las líneas de la entrada, salvo que las condiciones en el ámbito del comando indiquen lo contrario.

```
nereida:~/sed> cat b.without.sed
#example of succesive replacements
s/fish/cow/
s/cow/horse/
nereida:~/sed> cat b.test
fish
cow
nereida:~/sed> sed -f b.without.sed b.test
horse
horse
```

como se ve en el ejemplo, si un comando cambia la entrada, los siguientes comandos se aplican a la línea modificada (denominada *pattern space*).

Los comandos `sed` tienen el formato:

```
[direccion1][,direccion2][!]comando[argumentos]
```

1.1. Transferencia de Control

La orden `b` tiene la sintáxis:

```
[address1][,address2]b[label]
```

Transfiere el control a la etiqueta especificada. Si no se especifica la etiqueta, el control se transfiere al final del *script*, por lo que no se aplica ningún otro comando a la línea actual. Los dos siguientes ejemplos utilizan `b`:

```
$ cat b.sed
s/fish/cow/
b
s/cow/horse/
$ sed -f b.sed b.test
cow
cow
```

Utilizando *b* con una etiqueta:

```
$ cat blabel.sed
s/fish/cow/
b another
s/cow/horse/
:another
s/cow/cat/
$ sed -f blabel.sed b.test
cat
cat
```

1.2. Inserción de Texto

El uso de llaves {, } permite ejecutar los comandos en la lista entre llaves a las líneas seleccionadas. La llave que cierra debe estar en su propia línea aparte. Las llaves nos permiten, como se ve en el ejemplo, anidar selecciones y expresar condiciones del tipo “si esta entre estos dos patrones y además está entre estos otros dos ...”.

Los comandos *a* e *i* tienen una sintaxis parecida:

$[address]a\backslash$	$[address]i\backslash$
<i>text</i>	<i>text</i>

a añade (*i* inserta) el texto en cada línea que casa con la dirección especificada. El *text* no queda disponible en el “*pattern space*”, de manera que los subsiguientes comandos no le afectan. El siguiente ejemplo convierte un fichero *ascii* a *html*:

```
$ cat aandi.sed
1{
i\
<html>\
<head>\
<title>
p
i\
</title>\
</head>\
<body bgcolor=white>
}
$a\
</pre>\
</body>\
</html>
$ cat aandi.test
hello.world!
$ sed -f aandi.sed aandi.test
<html>
<head>
<title>
hello.world!
</title>
</head>
<body bgcolor=white>
hello.world!
</pre>
```

```
</body>
</html>
```

1.3. Trasferencia de Control Condicional

La sintaxis de la orden *t* es:

[address1][/,address2]t[label]

Si la sustitución tiene éxito, se bifurca a *label*, Si la etiqueta no se especifica, se salta al final del *script*. Consideremos el siguiente fichero de entrada:

```
$ cat t.test
name: "fulano de tal" address: "Leganitos,4" phone: "342255"
name: "fulano de cual"
name: "fulano de alli" address: "Legitos,4"
name: "zutano de tal" address: "Leg,8" phone: "342255"
```

Se asume que siempre que figura el teléfono se ha puesto la dirección y que si está la dirección se ha escrito el nombre. Se pretenden rellenar las líneas con un campo por defecto:

```
$ cat t.sed
/name:/{
s/. *name:[ ]*"."*[ ]*address:[ ]*"."*[ ]*. *phone:[ ]*".".*&/&/
t
s/. *name:[ ]*"."*[ ]*address:[ ]*"."*& phone: "????"/
t
s/. *name:[ ]*"."*& address: "???? ????"/ phone: "????"/
}
```

Esta es la llamada al script y la salida generada:

```
$ sed -f t.sed t.test
name: "fulano de tal" address: "Leganitos,4" phone: "342255"
name: "fulano de cual" address: "???? ????"/ phone: "????"/
name: "fulano de alli" address: "Legitos,4" phone: "????"/
name: "zutano de tal" address: "Leg,8" phone: "342255"
$
```

Ejemplo:

El fichero de entrada es un *folder* de pine en el que los mensajes recibidos tienen el formato:

...

```
01_NAME: XXX
02_SURNAME: XXX
03_TITLE: SISTEMAS
04_OtherTitle:
05_BIRTHDAY: XXX
06_BIRTHMONTH: XXX
07_BIRTHYEAR: XXX
08_ADDRESSFAMILY: XXX
09_ADDRESSACTUAL: XXX
10_POSTALCODE: XXX
11_EMAIL: XXX@csi.u11.es
12_TELEPHONE: XXX
13_FAX: XXX
14_LABGROUP: xxx
15_COMMENTS:
```

...

Se trata de escribir un script que produzca como salida los *emails* de los diferentes alumnos. Esta es una solución (suponemos que el *script* se llama con la opción `-n`):

```
#!/bin/sed -f
/^11_EMAIL:/{
s/11_EMAIL: *\([a-zA-Z0-9]*@[a-zA-Z0-9.]*\)\/\1/
t print
s/11_EMAIL: *\([a-zA-Z0-9]*\)\/\1@csi.ucll.es/
:print
p
}
```

Una característica no comentada y observada en algunos `sed`, incluyendo la versión Linux, es que, si existen varias sustituciones consecutivas antes de la sentencia de transferencia de control, basta con que una tenga éxito para que el salto se realice.

1.4. Rangos

Cuando se usan dos direcciones separadas por una coma, el rango que representan se extiende desde la primera línea que casa con el patrón hasta la siguiente línea que casa con el segundo patrón. El siguiente *script* muestra las tablas que aparecen en un fichero `LATEX`:

```
$ cat tables.sed
#Imprime las tablas en un fichero tex
/\begin{tabular}/,/end{tabular}/p
$ sed -n -f tables.sed *.tex
\begin{tabular}{c|c}
      [address]a$\backslash$ & [address]i$\backslash$\\
      text                  & text\\
\end{tabular}\\
\begin{tabular}{c|c}
      [address]a$\backslash$ & [address]i$\backslash$\\
      text                  & text\\
\end{tabular}\\
```

La siguiente selección para un rango comienza después de la última línea del rango previo; esto es, si el rango es `/A/,B/`, el primer conjunto de líneas que casa va desde la primera aparición de `/A/` hasta la siguiente aparición de `/B/`. Así, el número mínimo de líneas seleccionadas (si no es cero) es dos. Sólo en el caso en que la primera dirección es una expresión regular y la segunda un número de línea que viene antes de la línea que casa, es que sólo se selecciona una línea de emparejamiento. No hay por tanto solapes entre dos casamientos en un mismo rango de direcciones. El siguiente ejemplo ilustra la forma en la que `sed` interpreta los rangos de direcciones. Es una aproximación al problema de escribir los comentarios de un programa *C*.

```
$ cat scope.sed
#execute: sed -n -f scope.sed scope.test
/\/*/,/*\//p
```

Este *script* puede funcionar con programas cuyos comentarios (no anidados) empiecen y terminen en líneas diferentes, como el del ejemplo:

```
$ cat scope2.test
#file scope2.test
#include <stdio.h>
```



```

fact(int n) { /* recursive function
if(n == 0) return 1;
else(return n*fact(n-1));*/
}

void toto(id) {

main() {
    toto();
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
}

```

La ejecución del script selecciona los dos grupos de líneas que contienen comentarios:

```

$ sed -n -f scope.sed scope2.test
fact(int n) { /* recursive function
if(n == 0) return 1;
else(return n*fact(n-1));*/
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */

```

Sin embargo, el script fallará en este otro ejemplo:

```

$ cat scope.test
#include <stdio.h>

fact(int n) { /* recursive function */
if(n == 0) return 1;
else(return n*fact(n-1));
}

void toto(id) {

main() {
    toto();
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
}

```

La ejecución del programa da como resultado:

```

$ sed -n -f scope.sed scope.test
fact(int n) { /* recursive function */
if(n == 0) return 1;
else(return n*fact(n-1));
}

void toto(id) {

main() {
    toto();
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
}

```

1.5. Uso de la negación

El ejemplo muestra como convertir los caracteres españoles en un fichero \LaTeX escrito usando un teclado español a sus secuencias \LaTeX . puesto que en modo matemático los acentos tienen un significado distinto, se aprovecha la disponibilidad de las teclas españolas para simplificar la labor de tecleado. Obsérvese el uso de la exclamación ! para indicar la negación del rango de direcciones seleccionadas.

```
> cat sp2en.sed
/!begin{verbatim}/,/end{verbatim}/!{
/!begin{math}/,/end{math}/!{
s/á/\\'a/g
s/é/\\'e/g
s/í/\\'\{\\i\}/g
s/ó/\\'o/g
s/ú/\\'u/g
s/â/\\^a/g
s/ç/\\c{c}/g
s/ñ/\\~n/g
s/Á/\\'A/g
s/É/\\'E/g
s/Í/\\'I/g
s/Ó/\\'O/g
s/Ú/\\'U/g
s/Ñ/\\~N/g
s/¿/ > /g
s/¡/ < /g
}
}
/!begin{math}/,/end{math}/!{
s/â/\\hat{a}/g
s/á/\\acute{a}/g
s/à/\\grave{a}/g
s/ä/\\ddot{a}/g
s/ê/\\hat{e}/g
s/é/\\acute{e}/g
s/è/\\grave{e}/g
s/ë/\\ddot{e}/g
s/ô/\\hat{o}/g
s/ó/\\acute{o}/g
s/ò/\\grave{o}/g
s/ö/\\ddot{o}/g
s/û/\\hat{u}/g
s/ú/\\acute{u}/g
s/ù/\\grave{u}/g
s/ü/\\ddot{u}/g
}
```

Supongamos el siguiente fichero \LaTeX de entrada:

```
> cat sp2en.tex
\documentclass[11pt,a4paper,oneside,onecolumn]{article}
\usepackage{isolatin1}
\title{Stream Editor. Convirtiendo a Ingles en LaTeX}
\author{Casiano R. León \thanks{DEIOC Universidad de La Laguna}}
```

```

\begin{document}
\maketitle
Esto es un ejemplo de uso de sp2en.sed:
\begin{center}
áéíóú ÁÉÍÓÚ ñÑ ¿? ¡!\
\begin{math}
\hat{a}^{\acute{e}^2} \neq \hat{a}^{2\acute{e}}
\end{math}
\end{center}
comienza el verbatim\
\begin{listing}{1}
Lo que salga en verbatim depende de los packages
que hayan sido cargados: \a \e áéíóú ÁÉÍÓÚ ñÑ ¿? ¡!
\end{verbatim}
Termina el verbatim:\
\begin{center}
áéíóú ÁÉÍÓÚ ñÑ ¿? ¡!
\end{center}
\end{document}

```

Al ejecutar:

```
> sed -f sp2en.sed sp2en.tex
```

Obtenemos la salida:

```

\documentclass[11pt,a4paper,oneside,onecolumn]{article}
\usepackage{isolatin1}
\title{Stream Editor. Convirtiendo a Ingles en LaTeX}
\author{Casiano R. Le'on \thanks{DEIOC Universidad de La Laguna}}
\begin{document}
\maketitle
Esto es un ejemplo de uso de sp2en.sed:
\begin{center}
\ 'a\ 'e\ '{\i}\ 'o\ 'u \ 'A\ 'E\ 'I\ 'O\ 'U \ ~n\ ~N >? <!\
\begin{math}
\hat{a}^{\acute{e}^2} \neq \hat{a}^{2\acute{e}}
\end{math}
\end{center}
comienza el verbatim\
\begin{listing}{1}
Lo que salga en verbatim depende de los packages
que hayan sido cargados: \a \e áéíóú ÁÉÍÓÚ ñÑ ¿? ¡!
\end{verbatim}
Termina el verbatim:\
\begin{center}
\ 'a\ 'e\ '{\i}\ 'o\ 'u \ 'A\ 'E\ 'I\ 'O\ 'U \ ~n\ ~N >? <!
\end{center}
\end{document}

```

1.6. Siguiendo Línea: La orden n

Mediante la orden n podemos reemplazar el espacio de patrones actual por la siguiente línea. Si no se ha utilizado la opción `-n` en la ejecución de `sed`, el contenido del espacio de patrones se imprimirá antes de ser eliminada. El control pasa al comando siguiendo al comando n y no al primer

comando del *script*. Se trata, por tanto, de una orden que altera la conducta habitual de *sed*: Lo común es que, cuando se lee una nueva línea, se ejecute el primer comando del guión.

El siguiente ejemplo extrae los comentarios de un programa *C*. Se asume que, aunque los comentarios se pueden extender sobre varias líneas, no existe mas de un comentario por línea.

```
$ cat n.sed
# run it with: sed -n -f n.sed n.test
# write commented lines in a C program
#If current line matches /* ...
/\/*/{
# Comienzo de comentario, aseguremonos que la línea no casa con un cierre
# de comentario.
:loop
/\*\\!{
p
n
b loop
}
p
}
```

Supongamos el siguiente programa de prueba:

```
$ cat n.test
#include <stdio.h>

fact(int n) { /* recursive function */
if(n == 0) return 1;
else(return n*fact(n-1));
}

void toto(id) { /* This function */ /* is
                still empty */
}

main() {
    toto();
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
    toto();
    /* and here is
       a comment
       extended on 3 lines
    */
}
```

La salida al ejecutar el programa es la siguiente:

```
$ sed -n -f n.sed n.test
fact(int n) { /* recursive function */
void toto(id) { /* This function */ /* is
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
    /* and here is
       a comment
```

```

    extended on 3 lines
*/

```

Observe la desaparición de la línea “ still empty */” debido a la existencia de dos comentarios, uno de ellos inacabado en la línea anterior.

1.7. Manipulando tablas numéricas

El siguiente ejemplo intercambia la primera y ultima columnas de un fichero como el que sigue:

```

$ cat columns.test
11111 22222 33333 44444 55555
  111   22   33  4444  5555
11111 22222 33333 44444 55555
  11   22   3   444  5555

```

La parte mas complicada es preservar el sangrado. El truco reside, en parte, en el patrón central \2, que memoriza los blancos después de la primera columna y un sólo blanco antes de la última.

```

$ cat columns.sed
s/^\( *[0-9][0-9]*\)\( *[0-9] \)\( *[0-9][0-9]*\)$/\3\2\1/
$ sed -f columns.sed columns.test
55555 22222 33333 44444 11111
  5555   22   33  4444   111
55555 22222 33333 44444 11111
  5555   22   3   444   11

```

El siguiente ejemplo utiliza una opción del operador de sustitución que permite decidir que aparición del patrón deseamos sustituir. Así,

```

s/A/B/3

```

sustituirá la 3 aparición de A por B, obviando las otras.

El ejemplo selecciona la columna dos del fichero:

```

$ cat col2.sed
#extracts the second column
s/^ *//
s/\<[0-9][0-9]*\>//1
s/\<[0-9][0-9]*\>//2
s/\<[0-9][0-9]*\>//2
s/\<[0-9][0-9]*\>//2
s/ *$//
$ sed -f col2.sed columns.test
22222
  22
22222
  22

```

Mas general que el anterior, el siguiente ejemplo elimina un número de columnas arbitrario \$1 por la izquierda y otro número \$2 por la derecha. Para lograrlo, es necesario utilizar un guión para la shell que llama al correspondiente guión sed. Los parámetros son introducidos en el guión sed mediante el uso apropiado de las comillas dobles y simples:

```

> cat colbranch.sh
#!/bin/bash
sed -e '
s/^\( *[0-9]\+\)\{\ "$1" '\}\} //
s/\( *[0-9]\+\)\{\ "$2" '\}\} //
,

```

Veamos un ejemplo de uso:

```
> cat columns.test
11111 22222 33333 44444 55555
  111    22    33  4444  5555
11111 22222 33333 44444 55555
  11    22    3   444  5555
> colbranch.sh 2 1 < columns.test
33333 44444
   33  4444
33333 44444
   3   444
```

1.8. Traducción entre Tablas

El comando *y* realiza una traducción entre dos tablas de caracteres. Observa el ejemplo:

```
$ cat toupper.sed
y/aáéíóúäëïöübcdefghijklmnopqrstuvwxyzñ/AÁÉÍÓÜÄËÏÖÜBCDEFGHIJKLMNOPQRSTUVWXYZÑ/
$ cat sp2en.test
¡Coño! ¿Es plím el que hizo plúm?
```

Obtenemos así los contenidos del fichero en mayúsculas:

```
$ sed -f toupper.sed sp2en.test
¡COÑO! ¿ES PLÍM EL QUE HIZO PLÚM?
```

1.9. Del espacio de Patrones al de Mantenimiento

En *sed* se dispone, como en muchos otros editores de una zona de memoria a la que se puede enviar texto “cortado” o ‘copiado’ y desde la cual se puede recuperar el texto para insertarlo en la zona de trabajo (*pattern space*). En la jerga *sed* dicho espacio se conoce como *hold space*. El contenido del espacio de patrones (*pattern space*) puede moverse al espacio de mantenimiento (*hold space*) y recíprocamente:

<i>Hold</i>	h ó H	Copia o añade (append) los contenidos del <i>pattern space</i> al <i>hold space</i> .
<i>Get</i>	g ó G	Copia o añade los contenidos del <i>hold space</i> al <i>pattern space</i> .
<i>Exchange</i>	x	Intercambia los contenidos del <i>hold space</i> y el <i>pattern space</i>

Los comandos en minúsculas sobrescriben mientras que los que van en mayúsculas añaden. Así *h* copia los contenidos del *pattern space* en el *hold space*, borrando los contenidos previos que estuvieran en el *hold space*. Las orden *G* añade (paste) los contenidos del *hold space* al espacio de patrones actual (por el final). Las dos cadenas se enlazan a través de un retorno de carro.

Ejemplo Este guión intenta mediante una heurística poner la definiciones de funciones C al final del fichero, suponiendo que una definición de función comienza en la columna 1 y que se caracterizan mediante uno o dos identificadores seguidos de un paréntesis:

```
$ cat G.sed
/\<if>/s/./&/
t
/\<else>/s/./&/
t
#... lo mismo para las otras palabras clave
```

```

/^ [a-zA-Z] [a-zA-Z]* ([A-Z])*/H
t
/^ [a-zA-Z] [a-zA-Z]* * [a-zA-Z] [a-zA-Z]* (/H
${
G
}

```

Ejemplo de uso:

```

$ cat p.test
#include <stdio.h>

fact(int n) { /* recursive function */
if(n == 0) return 1;
else(return n*fact(n-1));
}

void toto(id) {

main() {
    toto();
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
}

```

Al ejecutar nuestro *script* obtenemos la salida:

```

$ sed -f G.sed p.test
#include <stdio.h>

fact(int n) { /* recursive function */
if(n == 0) return 1;
else(return n*fact(n-1));
}

void toto(id) {

main() {
    toto();
    printf("Hello world! factorial of 4 = %d\n",fact(4)); /* the
        comment takes two lines */
}

fact(int n) { /* recursive function */
void toto(id) {
main() {

```

Ejemplo El siguiente guión invierte los contenidos de un fichero:

```

> cat reverse.sed
${!
:loop
h
n

```

```
G
$!b loop
}
${
P
}
```

He aqui el resultado de una ejecución:

```
> cat reverse.test
one
two
three
> sed -n -f reverse.sed reverse.test
three
two
one
```

Ejemplo

Supuesto un fichero de entrada con formato: ...

```
NAME: xxx
SURNAME: xxx
TITLE: SISTEMAS
OtherTitle:
BIRTHDAY: xxx
BIRTHMONTH: xxx
BIRTHYEAR: xxx
ADDRESSFAMILY: xxx
ADDRESSACTUAL: xxx
POSTALCODE: xxx
EMAIL: XXXXXXXX@csi.ull.es
TELEPHONE: xxx
FAX: xxx
LABGROUP: xxx
COMMENTS:
```

...

Se pretenden extraer los apellidos y el nombre, concatenados con una coma. He aqui una posible solución:

```
#!/bin/sed -f
/^NAME:/ {
    s/^NAME://
    h
    n
    s/^SURNAME://
    G
    s/\n/,/
    y/áéíóúabcdefghijklmnñopqrstuvwxyz/ÁÉÍÓÚABCDEFGHIJKLMNÑOPQRSTUVWXYZ/
    P
}
```


1.10. La orden *N*

N añade la siguiente línea de entrada al espacio de patrones. Las dos líneas se separan mediante un retorno de carro. Después de su ejecución, el comando *N* pasa el control a los subsiguientes comandos en el *script*.

El siguiente ejemplo propuesto en [?] muestra una búsqueda y sustitución multilínea. Se trata de sustituir la cadena “Owner and Operator Guide” por “Installation Guide”, cualquiera que sea su posición en una línea o entre ellas. Los autores de [?] y [?] proponen la siguiente solución:

```
$ cat multiline2.sed
#Assume the pattern is in no more than two lines
s/Owner and Operator Guide/Installation Guide/g
/Owner/{
  N
  s/ *\n/ /g
  s/Owner *and *Operator *Guide/Installation Guide/g
}
```

Veamos primero los contenidos del fichero de prueba:

```
$ cat multiline.test
Dear Owner: Consult
Section 3.1 in the Owner and
Operator Guide for a description of the tape drives available for the Owner
of your system.
```

```
Consult Section 3.1 in the Owner and Operator
Guide for a description of the tape drives
available on your system.
```

```
Look in the Owner and Operator Guide, we mean the Owner
and Operator Guide shipped with your system.
```

```
Two manuals are provided including the Owner and
Operator Guide and the User Guide.
```

```
The Owner and Operator Guide is shipped with your system.
```

```
Look in the Owner
and Operator Guide shipped with your system.
```

```
The Owner
and
Operator
Guide is shipped with your system.
```

La ejecución del *script* da la siguiente salida:

```
$ sed -f multiline2.sed multiline.test
Dear Owner: Consult Section 3.1 in the Owner and
Operator Guide for a description of the tape drives available for the Owner
of your system.
```

```
Consult Section 3.1 in the Installation Guide for a description of the tape
drives
available on your system.
```

Look in the Installation Guide, we mean the Installation Guide shipped with your system.

Two manuals are provided including the Installation Guide and the User Guide.

The Installation Guide is shipped with your system.

Look in the Installation Guide shipped with your system.

The Owner and
Operator
Guide is shipped with your system.

Uno de los problemas, que aparece en el primer párrafo del ejemplo de prueba, es que la segunda línea leída debe ser reciclada para su uso en la siguiente búsqueda. El segundo fallo, que aparece en el último párrafo, es consecuencia de la limitación del *script* para trabajar con patrones partidos en más de dos líneas.

Consideremos esta otra solución:

```
$ cat multiline.sed
s/Owner and Operator Guide/Installation Guide/g
/Owner/{
:label
N
s/\n/ /g
s/Owner *and *Operator *Guide/Installation Guide/g
/Owner *$/b label
/Owner *and *$/b label
/Owner *and *Operator *$/b label
}
```

Este otro *script* hace que *sed* permanezca en un bucle mientras la línea adjuntada en segundo lugar contenga un prefijo estricto de la cadena buscada.

```
$sed -f multiline.sed multiline.test
Dear Owner: Consult Section 3.1 in the Installation Guide for \
a description of the tape drives available for the Owner of \
your system.
```

Consult Section 3.1 in the Installation Guide for a description of the tape drives available on your system.

Look in the Installation Guide, we mean the Installation Guide \
shipped with your system.

Two manuals are provided including the Installation Guide and the User Guide.

The Installation Guide is shipped with your system.

Look in the Installation Guide shipped with your system.

The Installation Guide is shipped with your system.

Un problema que aparece con esta aproximación es la presencia de líneas muy largas. Las líneas permanecen en el espacio de trabajo mientras terminen en un prefijo de la cadena buscada. Para que

la salida quepa en la hoja he tenido que partir las líneas del fichero de salida, lo que he indicado con los símbolos \. Considere esta modificación:

```
#!/bin/sed -f
s/Owner and Operator Guide/Installation Guide/g
/Owner/{
:label
N
s/Owner\([ \n]*\)and\([ \n]*\)Operator\([ \n]*\)Guide/Installation\1\2Guide\3/g
/Owner *$/b label
/Owner  *and *$/b label
/Owner  *and  *Operator *$/b label
}
```

Es indudable la ventaja de disponer de esta capacidad de búsqueda multilínea no puede realizarse con otras utilidades como *ex* o *vi*.

1.11. Suprimir: El Comando *D*

El comando *D* suprime la primera parte (hasta el retorno de carro empotrado) en un espacio de patrones multilínea y bifurca al primer comando en el *script*. El retorno de carro empotrado puede describirse mediante la secuencia de escape `\n`. En el caso en que el espacio de patrones quede vacío como consecuencia de la supresión, se lee una nueva línea.

El siguiente ejemplo compacta una secuencia de líneas vacías en una sólo línea vacía.

```
1 > cat N.sed
2 /^${/
3 N
4 /\n$/D
5 }
```

Si la línea es vacía se lee la línea siguiente. Si esta también es vacía el espacio de patrones contiene `^\n$`. La orden *D* deja en el espacio de trabajo una línea vacía y bifurca al comienzo del *script* (sin que se lea una nueva línea). Por tanto nada ha sido impreso, no se ejecuta el comando final *p* que actúa por defecto. Como el espacio de trabajo contiene `^$`, “casa” con el patrón especificado en línea 2 y se lee la siguiente línea. Si esta nueva línea es no vacía, no se ejecutará la orden *D* de la línea 4 y si que lo hará la orden por defecto final, imprimiéndose la línea vacía y la nueva línea no vacía.

Al ejecutar este “*script*” sobre un fichero conteniendo una secuencia de líneas en blanco:

```
> cat N.test
one empty

two empty lines

three empty lines
```

end of file

Se obtiene el resultado:

```
> sed -f N.sed N.test
one empty
```

two empty lines

three empty lines

end of file

Un buen ejercicio es intentar predecir la conducta de esta otra solución alternativa, en la que la supresión *D* es sustituida por la *d*:

```
/^${/{  
N  
/^\\n$/d  
}
```

¿Qué ocurrirá? ¿Es correcta esta segunda solución?

1.12. Búsqueda entre líneas

Este otro caso, también está tomado de [?] y [?]. Se trata de extender la capacidad de búsqueda de *grep*, de modo que el patrón pueda ser encontrado incluso si se encuentra diseminado entre a lo más dos líneas. El *script* presentado es una ligera variante del que aparece en [?] y escribe la(s) líneas que casan precedidas del número de la segunda línea.

```
$ cat phrase  
#!/bin/sh  
# phrase -- search for words across two lines.  
# Prints the line number  
# $1 = search string; remaining args = filenames  
search=$1  
shift  
for file  
do  
sed -n '  
  /"$search"/b final  
  N  
  h  
  s/.*\\n//  
  /"$search"/b final  
  g  
  s/ *\\n//  
  /"$search"/{  
    g  
    b final  
  }  
  g  
  D  
:final  
=  
p  
' $file  
done
```

Así, con el ejemplo “multiline.test” usado anteriormente obtenemos la siguiente salida:

```

$ phrase "Owner and Operator Guide" multiline.test
3
Section 3.1 in the Owner and
Operator Guide for a description of the tape drives available for the Owner
7
Consult Section 3.1 in the Owner and Operator
Guide for a description of the tape drives
10
Look in the Owner and Operator Guide, we mean the Owner
14
Two manuals are provided including the Owner and
Operator Guide and the User Guide.
16
The Owner and Operator Guide is shipped with your system.
19
Look in the Owner
and Operator Guide shipped with your system.

```

Primero se busca el patrón `/'"$search"'/` en la línea actual. Observe el habilidoso uso de las comillas simples y dobles para permitir la sustitución de la variable. La primera comilla simple cierra la comilla simple al final de la línea 10. Las comillas dobles hacen que la *shell* sustituya `$search` por su valor. La nueva comilla simple permite continuar el texto sin sustituciones.

Si se encuentra el patrón de búsqueda, imprimimos el número de línea (comando `=`) y la línea. Si no, leemos la siguiente línea, formando un patrón multilínea. Salvamos las dos líneas en el *hold space*. Entonces intentamos buscar el patrón `/'"$search"'/` en la línea que acaba de incorporarse. Es por eso que eliminamos del espacio de patrones la primera línea con la orden `s/ *\n/`. Si se encuentra, imprimimos y se repite el ciclo. Si no, recuperamos las dos líneas del *hold space* sustituimos el retorno de carro por un blanco y realizamos una nueva búsqueda. Si tiene éxito, se obtienen las dos líneas y se imprimen. En caso contrario, esto es, si el patrón no se ha encontrado en ninguna de las dos líneas, es necesario preservar la última para el siguiente ciclo. Por eso, se obtienen una vez mas las líneas del *hold space* y se suprime con `D` la primera de ellas. Dado que `D` devuelve el control al comienzo del *script*, la segunda línea no es eliminada. De todos modos, el *script* no es capaz de captar cuando un prefijo del patrón aparece al final de esta segunda línea, como muestra el ejemplo de prueba. En el primer párrafo el patrón se encuentra dos veces y sólo es encontrado una.

1.13. Seleccionando Items en un Registro Multilínea

El ejercicio resuelto aquí consiste en listar los alumnos que han seleccionado un determinado grupo de prácticas. Suponemos la organización del fichero de entrada descrita en la sección 1.9. El *script* recibe como primer argumento el nombre del fichero conteniendo la carpeta de correo (**pine**) asociada con la asignatura y como segundo argumento el grupo. Un primer *script* que no es descrito aquí, denominado `makefichas.sed` produce como salida el archivo con la estructura descrita en la sección 1.9. El segundo guión, denominado `grupo.sh` y que es el que nos ocupa, produce como salida los alumnos que pertenecen a ese grupo ed prácticas.

Estos son los contenidos del *script* grupo:

```
~/bin/makefichas.sed -n ~/mail/$1 | grupo.sh $2 | sort -u
```

Los contenidos del fichero `grupo.sh` son:

```

1  #!/bin/bash
2  search=$1
3  sed -n '
4  /^NAME:/ {
5      s/^NAME://

```

```

6    h
7    n
8    s/^SURNAME://
9    G
10   s/\n/,/
11   y/ÁÉÍÓÚáéíóúabcdefghijklmnopqrstuvwxyz/AEIOUAEIOUABCDEFGHIJKLMNOPQRSTUVWXYZ/
12   h
13 }
14 /^LABGROUP:/ {
15   y/ÁÉÍÓÚáéíóúabcdefghijklmnopqrstuvwxyz/AEIOUAEIOUABCDEFGHIJKLMNOPQRSTUVWXYZ/
16   s/'"$search"'/&/
17   t print
18   b
19 :print
20   g
21   p
22 }
23 '

```

De nuevo hacemos uso de las comillas simples y dobles en este ejemplo. Obsérvese como se protege el guión `sed` entre las líneas 3 y 16. En la línea 16 el cierre de la comilla simple y el uso de la doble comilla permite la actuación de la interpretación de la `shell`, sustituyendo `$search` que coincide con el parámetro pasado en la llamada como `$2`. La siguiente comilla simple en esa línea permite la protección del resto del guión.

Capítulo 2

Expresiones Regulares en Flex

Un lenguaje regular es aquel que puede ser descrito mediante expresiones regulares como las que se utilizan en `ex`, `vi`, `sed`, `perl` y en tantas otras utilidades UNIX. Dado un lenguaje regular, un analizador léxico es un programa capaz de reconocer las entradas que pertenecen a dicho lenguaje y realizar las acciones semánticas que se hayan asociado con los estados de aceptación. Un generador de analizadores léxicos es una herramienta que facilita la construcción de un analizador léxico. Un generador de analizadores léxicos parte, por tanto, de un lenguaje adecuado para la descripción de lenguajes regulares (y de su semántica) y produce como salida una función (en C, por ejemplo) que materializa el correspondiente analizador léxico. La mayor parte de los generadores producen a partir del conjunto de expresiones regulares los correspondientes tablas de los autómatas finitos deterministas. Utilizando dichas tablas y un algoritmo de simulación genérico del autómata finito determinista se obtiene el analizador léxico. Una vez obtenido el estado de aceptación a partir de la entrada es posible, mediante una sentencia `switch` ejecutar la acción semántica asociada con la correspondiente expresión regular.

2.1. Estructura de un programa LEX

Estructura de un programa

LEX y FLEX son ejemplos de generadores léxicos. Flex lee desde la entrada estándar si no se especifica explícitamente un fichero de entrada. El fichero de entrada `reglen.1` (se suele usar el tipo 1) debe tener la forma:

```
%{  
declaration C1  
.  
.  
.  
  
declaration CM  
%}  
macro_name1 regular_definition1  
.  
.  
.  
  
macro_nameR regular_definitionR  
  
%x exclusive_state  
%s inclusive_state  
%%
```

```

regular_expression1 { action1(); }
.
.
.

regular_expressionN { actionN(); }

%%
support_routine1() {
}
.
.
.

support_routineS() {
}

```

Como vemos, un programa LEX consta de 3 secciones, separadas por `%%`. La primera sección se denomina *sección de definiciones*, la segunda *sección de reglas* y la tercera *sección de código*. La primera y la última son opcionales, así el programa legal LEX mas simple es:

```
%%
```

que genera un analizador que copia su entrada en `stdout`.

Compilación

Una vez compilado el fichero de entrada `regleng.l` mediante la correspondiente orden:

```
flex reglen.l
```

obtenemos un fichero denominado `lex.yy.c`. Este fichero contiene la rutina `yylex()` que realiza el análisis léxico del lenguaje descrito en `regleng.l`. Supuesto que una de las `support_routines` es una función `main()` que llama a la función `yylex()`, podemos compilar el fichero generado con un compilador C para obtener un ejecutable `a.out`:

```
cc lex.yy.c -lfl
```

La inclusión de la opción `-fl` enlaza con la librería de `flex`, que contiene dos funciones: `main` y `yywrap()`.

Ejecución

Cuando ejecutamos el programa `a.out`, la función `yylex()` analiza las entradas, buscando la secuencia mas larga que casa con alguna de las expresiones regulares (`regular_expressionK`) y ejecuta la correspondiente acción (`actionK()`). Si no se encuentra ningun emparejamiento se ejecuta la *regla por defecto*, que es:

```
(.|\n) { printf("%s",yytext); }
```

Si encuentran dos expresiones regulares con las que la cadena mas larga casa, elige la que figura primera en el programa `lex`. Una vez que `yylex()` ha encontrado el *token*, esto es, el patrón que casa con la cadena mas larga, dicha cadena queda disponible a través del puntero global `yytext`, y su longitud queda en la variable entera global `yylen`.

Una vez que se ha ejecutado la correspondiente acción, `yylex()` continúa con el resto de la entrada, buscando por subsiguientes emparejamientos. Asi continúa hasta encontrar un `end of file`, en cuyo caso termina, retornando un cero o bien hasta que una de las acciones explicitamente ejecuta una sentencia `return`.

Sección de definiciones

La primera sección contiene, si las hubiera, las definiciones regulares y las declaraciones de los estados de arranque.

Las definiciones tiene la forma:

name regular_definition

donde *name* puede ser descrito mediante la expresión regular:

`[a-zA-Z_][a-zA-Z_0-9-]*`

La *regular_definition* comienza en el primer carácter no blanco que sigue a *name* y termina al final de la línea. La definición es una expresión regular extendida. Las subsiguientes definiciones pueden “llamar” a la macro `{name}` escribiéndola entre llaves. La macro se expande entonces a (*regular_definition*) en *flex* y a *regular_definition* en *lex*.

El código entre los delimitadores `%{` y `%}` se copia verbatim al fichero de salida, situándose en la parte de declaraciones globales. Los delimitadores deben aparecer (sólos) al comienzo de la línea.

El Lenguaje de las Expresiones Regulares Flex

La sintáxis que puede utilizarse para la descripción de las expresiones regulares es la que se conoce como “extendida”:

- `x` Casa con 'x'
- `.` Cualquier carácter, excepto `\n`.
- `[xyz]` Una “clase”; en este caso una de las letras `x`, `y`, `z`
- `[abj-oZ]` Una “clase” con un rango; casa con `a`, `b`, cualquier letra desde `j` hasta `o`, o una `Z`
- `[^A-Z]` Una “Clase complementada” esto es, todos los caracteres que no están en la clase. Cualquier carácter, excepto las letras mayúsculas. Obsérvese que el retorno de carro `\n` casa con esta expresion. Así es posible que, ¡un patrón como `[^"]+` pueda casar con todo el fichero!.
- `[^A-Z\n]` Cualquier carácter, excepto las letras mayúsculas o un `\n`.
- `[[:alnum:]]` Casa con cualquier caracter alfanumérico. Aquí `[[:alnum:]]` se refiere a una de las clases predefinidas. Las otras clases son: `[[:alpha:]]` `[[:blank:]]` `[[:cntrl:]]` `[[:digit:]]` `[[:graph:]]` `[[:lower:]]` `[[:print:]]` `[[:punct:]]` `[[:space:]]` `[[:upper:]]` `[[:xdigit:]]`. Estas clases designan el mismo conjunto de caracteres que la correspondiente función C `isXXXX`.
- `r*` Cero o mas `r`.
- `r+` Una o mas `r`.
- `r?` Cero o una `r`.
- `r{2,5}` Entre 2 y 5 `r`.
- `r{2,}` 2 o mas `r`. `r{4}` Exactamente 4 `r`.
- `{macro_name}` La expansión de `macro_name` por su *regular_definition*
- `"[xyz]"` Exactamente la cadena: `[xyz]`
- `\X` Si `X` is an `a`, `b`, `f`, `n`, `r`, `t`, o `v`, entonces, la interpretación ANSI-C de `\x`. En cualquier otro caso `X`.
- `\0` El carácter NUL (ASCII 0).
- `\123` El carácter cuyo código octal es 123.
- `\x2a` El carácter cuyo código hexadecimal es 2a.
- `(r)` Los paréntesis son utilizados para cambiar la precedencia.

- **rs** Concatenation
- **r|s** Casa con **r** o **s**
- **r/s** Una **r** pero sólo si va seguida de una **s**. El texto casado con **s** se incluye a la hora de decidir cual es el emparejamiento mas largo, pero se devuelve a la entrada cuando se ejecuta la acción. La acción sólo ve el texto asociado con **r**. Este tipo de patrón se denomina *trailing context* o *lookahead* positivo.
- **^r** Casa con **r**, al comienzo de una línea. Un **^** que no aparece al comienzo de la línea o un **\$** que no aparece al final de la línea, pierde su naturaleza de “ancla” y es tratado como un carácter ordinario. Asi: `foo|(bar$)` se empareja con `bar$`. Si lo que se quería es la otra interpretación, es posible escribir `foo|(bar\n)`, o bien:

```
foo |
bar$  { /* action */ }
```

- **r\$** Casa con **r**, al final de una línea. Este es también un operador de *trailing context*. Una regla no puede tener mas de un operador de *trailing context*. Por ejemplo, la expresión `foo/bar$` es incorrecta.
- **<s>r** Casa con **r**, pero sólo si se está en el estado **s**.
- **<s1,s2,s3>r** Idem, si se esta en alguno de los estados **s1**, **s2**, or **s3**
- **<*>r** Casa con **r** cualquiera que sea el estado, incluso si este es exclusivo.
- **<<EOF>>** Un final de fichero.
- **<s1,s2><<EOF>>** Un final de fichero, si los estados son **s1** o **s2**

Los operadores han sido listados en orden de precedencia, de la mas alta a la mas baja. Por ejemplo `foo|bar+` es lo mismo que `(foo)|(ba(r)+)`.

Las Acciones Semánticas

Cada patrón regular tiene su correspondiente acción asociada. El patrón termina en el primer espacio en blanco (sin contar aquellos que están entre comillas dobles o prefijados de secuencias de escape). Si la acción comienza con `{`, entonces se puede extender a través de multiples líneas, hasta la correspondiente `}`. El programa `flex` no hace un análisis del código C dentro de la acción. Existen tres directivas que pueden insertarse dentro de las acciones: `BEGIN`, `ECHO` y `REJECT`. Su uso se muestra en los subsiguientes ejemplos.

La sección de código se copia verbatim en `lex.yy.c`. Es utilizada para proveer las funciones de apoyo que se requieran para la descripción de las acciones asociadas con los patrones que parecen en la sección de reglas.

2.2. Versión Utilizada

Todos los ejemplos que aparecen en este documento fueron preparados con la versión 2.5.4 de `flex` en un entorno Linux

```
$ uname -a
Linux nereida.deioc.ull.es 2.2.12-20 #10 Mon May 8 19:40:16 WEST 2000 i686 unknown
$ flex --version
flex version 2.5.4
```

y con la versión 2.5.2 en un entorno Solaris

```
> uname -a
SunOS fonil 5.7 Generic_106541-04 sun4u sparc SUNW,Ultra-5_10
> flex --version
flex version 2.5.2
```

2.3. Espacios en blanco dentro de la expresión regular

La expresión regular va desde el comienzo de la línea hasta el primer espacio en blanco no escapado. Todos los espacios en blanco que formen parte de la expresión regular deben ser escapados o protegidos entre comillas. Así, el siguiente programa produce un error en tiempo de compilación C:

```
> cat spaces.l
%%
one two { printf("spaces\n"; }
%%
nereida:~/public_html/regexpr/lex/src> flex spaces.l
nereida:~/public_html/regexpr/lex/src> gcc lex.yy.c
spaces.l: In function 'yylex':
spaces.l:2: 'two' undeclared (first use in this function)
spaces.l:2: (Each undeclared identifier is reported only once
spaces.l:2: for each function it appears in.)
spaces.l:2: parse error before '{'
spaces.l:4: case label not within a switch statement
lex.yy.c:632: case label not within a switch statement
lex.yy.c:635: case label not within a switch statement
lex.yy.c:757: default label not within a switch statement
lex.yy.c: At top level:
lex.yy.c:762: parse error before '}'
```

Deberíamos escapar el blanco entre `one` y `two` o bien proteger la cadena poniéndola entre comillas: `"one two"`.

2.4. Ejemplo Simple

Este primer ejemplo sustituye las apariciones de la palabra *username* por el *login* del usuario:

```
$ cat subst.l
%option main
%{
#include <unistd.h>
%}
%%
username    printf( "%s",  getlogin());
%%
$ flex -osubst.c subst.l
$ gcc -o subst subst.c
$ subst
Dear username:
Dear pl:
```

He presionado CTRL-D para finalizar la entrada.

Observe el uso de la opción `%option main` en el fichero *subst.l* para hacer que *flex* genere una función *main*. También merece especial atención el uso de la opción `-osubst` para cambiar el nombre del fichero de salida, que por defecto será *lex.yy.c*.

2.5. Suprimir

Al igual que en *sed* y *awk*, es muy sencillo suprimir las apariciones de una expresión regular.

```
$ cat delete.l
/* delete all entries of zap me */
%%
"zap me"
$ flex delete.l ; gcc lex.yy.c -lfl; a.out
this is zap me a first zap me phrase
this is  a first  phrase
```

2.6. Declaración de yytext

En la sección de definiciones es posible utilizar las directivas `%pointer` o `%array`. Estas directivas hacen que `yytext` se declare como un puntero o un *array* respectivamente. La opción por defecto es declararlo como un puntero, salvo que se haya usado la opción `-l` en la línea de comandos, para garantizar una mayor compatibilidad con LEX. Sin embargo, y aunque la opción `%pointer` es la mas eficiente (el análisis es mas rápido y se evitan los *buffer overflow*), limita la posible manipulación de *yytext* y de las llamadas a `unput()`.

```
$ cat yytextp.l
%%
hello {
    strcat(yytext, " world");
    printf("\n%d: %s\n",strlen(yytext),yytext);
}
$ flex yytextp.l ; gcc lex.yy.c -lfl ; a.out
hello

11: hello world

fatal flex scanner internal error--end of buffer missed
```

Este error no aparece si se utiliza la opción `%array`:

```
$ cat yytext.l
%array
%%
hello {
    strcat(yytext, " world");
    printf("\n%d: %s\n",strlen(yytext),yytext);
}
$ flex yytext.l; gcc lex.yy.c -lfl; a.out
hello

11: hello world
```

Además, algunos programs LEX modifican directamente `yytext`, utilizando la declaración:

```
extern char yytext[]
```

que es incompatible con la directiva `%pointer` (pero correcta con `%array`). La directiva `%array` define `yytext` como un *array* de tamaño `YYLMAX`. Si deseamos trabajar con un mayor tamaño, basta con redefinir `YYLMAX`.

2.7. Declaración de `yylex()`

Por defecto la función `yylex()` que realiza el análisis léxico es declarada como `int yylex()`. Es posible cambiar la declaración por defecto utilizando la macro `YY_DECL`. En el siguiente ejemplo la definición:

```
#define YY_DECL char *scanner(int *numcount, int *idcount)
```

hace que la rutina de análisis léxico pase a llamarse `scanner` y tenga dos parametros de entrada, retornando un valor de tipo `char *`.

```
$ cat decl.1
%{
#define YY_DECL char *scanner(int *numcount, int *idcount)
%}

num [0-9]+
id [a-z]+
%%
{num} {(*numcount)++;}
halt {return ((char *) strdup(yytext));}
{id} {(*idcount)++;}
%%
main() {
    int a,b;
    char *t;

    a = 0; b = 0;
    t = scanner(&a, &b);
    printf("numcount = %d, idcount = %d, yytext = %s\n",a,b,t);
    t = scanner(&a, &b);
    printf("numcount = %d, idcount = %d, yytext = %s\n",a,b,t);
}

int yywrap() {
    return 1;
}
```

La ejecución del programa anterior produce la siguiente salida:

```
$ decl
a b 1 2 3 halt
    numcount = 3, idcount = 2, yytext = halt

e 4 5 f

numcount = 5, idcount = 4, yytext = (null)
$ decl
a b 1 2 3 halt
    numcount = 3, idcount = 2, yytext = halt

e 4 f 5 halt
    numcount = 5, idcount = 4, yytext = halt
```

2.8. yywrap()

Cuando el analizador léxico alcanza el final del fichero, el comportamiento en las subsiguientes llamadas a *yylex* resulta indefinido. En el momento en que *yylex()* alcanza el final del fichero llama a la función *yywrap*, la cual retorna un valor de 0 o 1 según haya mas entrada o no. Si el valor es 0, la función *yylex* asume que la propia *yywrap* se ha encargado de abrir el nuevo fichero y asignárselo a *yyin* . Otra manera de continuar es haciendo uso de la función *yyrestart(FILE *file)*. El siguiente ejemplo cuenta el número de líneas, palabras y caracteres en una lista de ficheros proporcionados como entrada.

```
%{
unsigned long charCount = 0, wordCount = 0, lineCount = 0;
}%

word [^ \t\n]+
eol \n

%%
{word} { wordCount++; charCount += yyleng; }
{eol} { charCount++; lineCount++; }
. charCount++;

%%

char **fileList;
unsigned nFiles;
unsigned currentFile = 0;
unsigned long totalCC = 0;
unsigned long totalWC = 0;
unsigned long totalLC = 0;

main ( int argc, char **argv) {
    FILE *file;

    fileList = argv + 1; nFiles = argc - 1;

    if (nFiles == 0) {
        fprintf(stderr, "Usage is:\n%s file1 file2 file3 ...\n", argv[0]);
        exit(1);
    }
    file = fopen (fileList[0], "r");
    if (!file) {
        fprintf (stderr, "could not open %s\n", argv[1]);
        exit (1);
    }
    currentFile = 1; yyrestart(file);
    yylex ();
    printf ("%8lu %8lu %8lu %s\n", lineCount, wordCount,
        charCount, fileList[currentFile - 1]);
    if (argc > 2) {
        totalCC += charCount; totalWC += wordCount; totalLC += lineCount;
        printf ("%8lu %8lu %8lu total\n", totalLC, totalWC, totalCC);
    }
    return 0;
}
```

```

}

int yywrap () {
    FILE *file;

    if (currentFile < nFiles) {
        printf ("%8lu %8lu %8lu %s\n", lineCount, wordCount,
            charCount, fileList[currentFile - 1]);
        totalCC += charCount; totalWC += wordCount; totalLC += lineCount;
        charCount = wordCount = lineCount = 0;
        fclose (yyin);

        while (fileList[currentFile] != (char *) 0) {
            file = fopen (fileList[currentFile++], "r");
            if (file != NULL) { yyrestart(file); break; }
        }
        fprintf (stderr, "could not open %s\n", fileList[currentFile - 1]);
        return (file ? 0 : 1);
    }
    return 1;
}

```

La figura muestra el proceso de compilación y la ejecución:

```

$ flex countlwc.l;gcc lex.yy.c; a.out *.l
    58      179      1067 ape-05.l
    88      249      1759 countlwc.l
    11       21       126 magic.l
     9       17       139 mgrep.l
     9       16       135 mlg.l
     5       15       181 ml.l
     7       12        87 subst.l
   187      509     3494 total

```

La diferencia esencial entre asignar *yyin* o llamar a la función *yyrestart* es que esta última puede ser utilizada para conmutar entre ficheros en medio de un análisis léxico. El funcionamiento del programa anterior no se modifica si se se intercambian asignaciones a *yyin* (*yyin = file*) y llamadas a *yyrestart(file)*.

2.9. unput()

La función *unput(c)* coloca el carácter *c* en el flujo de entrada, de manera que será el primer carácter leído en próxima ocasión.

```

$ cat unput2.l
%array
%%
[a-z] {unput(toupper(yytext[0]));}
[A-Z] ECHO;
%%
$ flex unput2.l ; gcc lex.yy.c -lfl;a.out
abcd
ABCD

```

Un problema importante con *unput* es que, cuando se utiliza la opción `%pointer`, las llamadas a *unput* destruyen los contenidos de *yytext*. Es por eso que, en el siguiente ejemplo se hace una copia de *yytext*. La otra alternativa es, por supuesto, usar la opción `%array`.

```
$ cat unput.l
%%
[0-9]+ {
    int i;
    char *yycopy = (char *) strdup(yytext);

    unput(')');
    for(i=strlen(yycopy)-1; i>=0; --i)
        unput(yycopy[i]);
    unput('(');
    free(yycopy);
}
\[([0-9]+\\) printf("Num inside parenthesis: %s\n",yytext);
.\|n
$ flex unput.l ; gcc lex.yy.c -lfl ; a.out
32
Num inside parenthesis: (32)
(43)
Num inside parenthesis: (43)
```

2.10. input()

La función *input()* lee desde el flujo de entrada el siguiente carácter. Normalmente la utilizaremos si queremos tomar “personalmente el control” del análisis. El ejemplo permite “engullir” los comentarios (no anidados):

```
$ cat input.l
%%
"/*" {
    int c;
    for(;;) {
        while ((c=input()) != '*' && c != EOF)
            ;
        if (c == '*') {
            while ((c = input()) == '*')
                ;
            if (c == '/') break;
        }
        if (c == EOF) {
            fprintf(stderr,"Error: EOF in comment");
            yyterminate();
        }
    }
}
```

La función *yyterminate()* termina la rutina de análisis léxico y devuelve un cero indicándole a la rutina que llama que todo se ha acabado. Por defecto, *yyterminate()* es llamada cuando se encuentra un final de fichero. Es una macro y puede ser redefinida.

```
$ flex input.l ; gcc lex.yy.c -lfl ; a.out
hello /* world */
```



```
hello
unfinished /* comment
unfinished Error: EOF in comment
```

He presionado CTRL-D después de entrar la palabra *comment*.

2.11. REJECT

La directiva **REJECT** le indica al analizador que proceda con la siguiente regla que casa con un prefijo de la entrada. Como es habitual en *flex*, se elige la siguiente regla que casa con la cadena mas larga. Consideremos el siguiente ejemplo:

```
$ cat reject.l
%%
a      |
ab     |
abc    |
abcd ECHO; REJECT; printf("Never seen\n");
.| \n
```

La salida es:

```
$ gcc lex.yy.c -lfl;a.out
abcd
abcdabcaba
```

Observe que **REJECT** supone un cambio en el flujo de control: El código que figura después de **REJECT** no es ejecutado.

2.12. yymore()

La función `yymore()` hace que, en vez de vaciar *yytext* para el siguiente *matching*, el valor actual se mantenga, concatenando el valor actual de *yytext* con el siguiente:

```
$ cat yymore.l
%%
mega- ECHO; yymore();
kludge ECHO;

$ flex yymore.l ; gcc lex.yy.c -lfl ; a.out
mega-kludge
mega-mega-kludge
```

La variable `yylen` no debería ser modificada si se hace uso de la función `yymore()`.

2.13. yless()

La función `yless(n)` permite retrasar el puntero de lectura de manera que apunta al carácter *n* de *yytext*. Veamos un ejemplo:

```
$ cat yless.l
%%
foobar ECHO; yless(4);
[a-z]+ ECHO;
```

```
$ flex yyless.1; gcc lex.yy.c -lfl; a.out
foobar
foobarar
```

Veamos un ejemplo mas “real”. supongamos que tenemos que reconocer las cadenas entre comillas dobles, pero que pueden aparecer en las mismas secuencias de escape `\`. La estrategia general del algoritmo es utilizar la expresión regular `\("[^"]+"\)` y examinar si los dos últimos caracteres en `yytext` son `\`. En tal caso, se concatena la cadena actual (sin la `"` final) como prefijo para el próximo emparejamiento (utilizando `yymore`). La eliminación de la `"` se hace a través de la ejecución de `yyless(yytext-1)`, que al mismo tiempo garantiza que el próximo emparejamiento tendrá lugar con este mismo patrón `\("[^"]+"\)`.

```
$ cat quotes.1
%%
\[^[^"]+"\) {
    printf("Processing string. %d: %s\n",yytext,yytext);
    if (yytext[yytext-2] =='\') {
        yyless(yytext-1); /* so that it will match next time */
        yymore(); /* concatenate with current yytext */
        printf("After yyless. %d: %s\n",yytext,yytext);
    } else {
        printf("Finished. The string is: %s\n",yytext);
    }
}
```

El ejemplo no puede entenderse si no se tiene en cuenta que `yyless(yytext-1)` actualiza los valores de `yytext` y `yytext`, como muestra la salida.

¿Qué ocurre si intercambiamos las posiciones de `yymore()` e `yyless(yytext-1)` en el código? ¿Cambiaría la salida? La respuesta es que no. Parece que la concatenación se hace con el valor final de `yytext` y no con el valor que este tenía en el momento de la llamada a `yymore`.

Otra observación a tener en cuenta es que `yyless()` es una macro y que, por tanto, sólo puede ser utilizada dentro del fichero `lex` y no en otros ficheros fuentes.

En general, el uso de estas funciones nos puede resolver el problema de reconocer límites que de otra forma serían difíciles de expresar con una expresión regular.

```
$ flex quotes.1 ; gcc lex.yy.c -lfl ; a.out
"Hello \"Peter\", nice to meet you"
Processing string. 9: "Hello \"
After yyless. 8: "Hello \"
Processing string. 16: "Hello \"Peter\"
After yyless. 15: "Hello \"Peter\"
Processing string. 35: "Hello \"Peter\", nice to meet you"
Finished. The string is: "Hello \"Peter\", nice to meet you"
```

2.14. Estados

Las expresiones regulares pueden ser prefijadas mediante *estados*. Los estados o condiciones de arranque, se denotan mediante un identificador entre ángulos y se declaran en la parte de las definiciones. Las declaraciones se hacen mediante `%s` para los estados “inclusivos” o bien `%x` para los estados “exclusivos”, seguidos de los nombres de los estados. No pueden haber caracteres en blanco antes de la declaración. Un *estado* se activa mediante la acción `BEGIN estado`. A partir de ese momento, las reglas que esten prefijadas con el estado pasan a estar activas. En el caso de que el estado sea inclusivo, las reglas no prefijadas también permanecen activas. Los estados exclusivos son especialmente útiles para especificar “sub analizadores” que analizan porciones de la entrada cuya estructura “sintáctica” es diferente de la del resto de la entrada.

El ejemplo “absorbe” los comentarios, conservando el numero de líneas del fichero en la variable `linenum`

```
$ cat comments.1
%option noyywrap
%{
    int linenum = 0;
}%
%x comment
%%

"/" BEGIN(comment); printf("comment=%d, YY_START = %d, YYSTATE = %d",comment,YY_START,YYSTATE);
<comment>[^\n]* /* eat anything that is not a star * /
<comment>"*"+[^\n]* /* eat up starts not followed by / */
<comment>\n ++linenum; /* update number of lines */
<comment>"*"+"/" BEGIN(0);

\n ECHO; linenum++;
. ECHO;
%%
main() {
    yylex();
    printf("\n%d lines\n",linenum);
}
```

La opción `noyywrap` hace que `yylex()` no llame a la función `yywrap()` al final del fichero y que asuma que no hay mas entrada por procesar.

Los estados se traducen por enteros, pudiendo ser manipulados como tales. La macro `INITIAL` puede utilizarse para referirse al estado 0. Las macros `YY_START` y `YYSTATE` contienen el valor del estado actual.

```
$ flex comments.1 ; gcc lex.yy.c ; a.out < hello.c
main() <%
int a<:1:>; comment=1, YY_START = 1, YYSTATE = 1
    a<:0:> = 4; comment=1, YY_START = 1, YYSTATE = 1
    printf("hello world! a(0) is %d\n",a<:0:>);
%>
```

```
6 lines
$ cat hello.c
main() <%
int a<:1:>; /* a comment */
    a<:0:> = 4; /* a comment in
                two lines */
    printf("hello world! a(0) is %d\n",a<:0:>);
%>
```

En *flex* es posible asociar un ámbito con los estados o condiciones iniciales. Basta con colocar entre llaves las parejas *patrón acción* gobernadas por ese estado. El siguiente ejemplo procesa las cadenas *C*:

```
$ cat ststring.1
%option main
%x str
%{
```

```

#define MAX_STR_CONST 256

char string_buffer[MAX_STR_CONST];
char *string_buf_ptr;
%}

%%
\" string_buf_ptr = string_buffer; BEGIN(str);
<str>{
\"          {BEGIN (INITIAL); *string_buf_ptr = '\\0'; printf(\"%s\",string_buffer); }
\\n          { printf(\"Error: non terminated string\\n\"); exit(1); }
\\[0-7]{1,3} { int result; /* octal escape sequence */
              (void) sscanf(yytext+1,\"%o",&result);
              if (result > 0xff) {printf(\"Error: constant out of bounds\\n\"); exit(2); }
              *string_buf_ptr++ = result;
            }
\\[0-9]+      { printf(\"Error: bad escape sequence\\n\"); exit(2); }
\\\\n          { *string_buf_ptr++ = '\\n'; }
\\\\t          { *string_buf_ptr++ = '\\t'; }
\\\\b          { *string_buf_ptr++ = '\\b'; }
\\\\r          { *string_buf_ptr++ = '\\r'; }
\\\\f          { *string_buf_ptr++ = '\\f'; }
\\\\(.|\\n)      { *string_buf_ptr++ = yytext[1]; }
[^\\\\n\"']+    { char *yptr = yytext; while(*yptr) *string_buf_ptr++ = *yptr++; }
}
(.|\\n)
%%
$ flex ststring.l ; gcc lex.yy.c ; a.out < hello.c
    hello
world! a(0) is %d
$ cat hello.c
main() <%
int a<:1:>; /* a comment */
    a<:0:> = 4; /* a comment in
                two lines */
    printf(\"\\thell\\157\\nworld! a(0) is %d\\n\",a<:0:>);
%>

```

Observe la conducta del programa ante las siguientes entradas:

- Entrada:

```
"hello \
dolly"
```

¿Cuál será la salida? ¿Que patrón del programa anterior es el que casa aqui?

- Entrada: "hello\\ndolly". ¿Cuál será la salida? ¿Que patrón del programa anterior es el que casa aqui?

- |


```
"hello"
```

Donde hay un retorno del carro después de hello. ¿Cuál será la salida?

2.15. La pila de estados

Mediante el uso de la opción

`%option stack`

tendremos acceso a una pila de estados y a tres rutinas para manipularla:

- `void yy_push_state(int new_state)`
Empuja el estado actual y bifurca a `new_state`.
- `void yy_pop_state()`
Saca el estado en el *top* de la pila y bifurca a el mismo.
- `int yy_top_state()`
Nos devuelve el estado en el *top* de la pila, sin alterar los contenidos de la misma.

2.15.1. Ejemplo

El siguiente programa `flex` utiliza las funciones de la pila de estados para reconocer el lenguaje (no regular) $\{a^n b^n / n \in N\}$

```
%option main
%option noyywrap
%option stack
%{
#include <stdio.h>
#include <stdlib.h>
%}
%x estado_a
%%
^a { yy_push_state(estado_a); }
<estado_a>{
a      { yy_push_state(estado_a); }
b      { yy_pop_state(); }
b[~b\n]+ { printf ("Error\n");
          while (YYSTATE != INITIAL)
              yy_pop_state();
          while (input() != '\n') ;
      }
(.|\n) { printf ("Error\n");
          while (YYSTATE != INITIAL)
              yy_pop_state();
          while (input() != '\n') ;
      }
}
.      { printf ("Error\n");
          while (input() != '\n') ;
      }
\n     { printf("Aceptar\n"); }
%%
```

2.16. Final de Fichero

El patrón <<EOF>> permite asociar acciones que se deban ejecutar cuando se ha encontrado un *end of file* y la macro `yywrap()` ha devuelto un valor no nulo.

Cualquiera que sea, la acción asociada deberá de optar por una de estas cuatro alternativas:

- Asignar `yyin` a un nuevo fichero de entrada
- Ejecutar `return`
- Ejecutar `yyterminate()` (véase la sección 2.10)
- Cambiar de *buffer* de entrada utilizando la función `yy_switch_buffer` (véase la sección 2.21).

El patrón <<EOF>> no puede usarse con otras expresiones regulares. Sin embargo, es correcto prefijarlo con estados.

Si <<EOF>> aparece sin condiciones de arranque, la regla se aplica a todos los estados que no tienen una regla <<EOF>> específica. Si lo que se quiere es que la regla se restrinja al ámbito del estado inicial se deberá escribir:

```
<INITIAL><<EOF>>
```

Sigue un programa que reconoce los comentarios anidados en C. Para detectar comentarios incabados usaremos <<EOF>>.

```
%option stack
%x comment
%%
"/*" { yy_push_state(comment); }
(.|\n) ECHO;
<comment>"/*" { yy_push_state(comment); }
<comment>"/" { yy_pop_state(); }
<comment>(.|\n) ;
<comment><<EOF>> { fprintf(stderr,"Error\n"); exit(1); }
%%

$ cat hello.c
main() {
int a[1]; /* a /* nested comment */. */
    a[0] = 4; /* a /* nested comment in
                /* two */ lines */ *****/
}
$ flex nestedcom.l ; gcc lex.yy.c -lfl ; a.out < hello.c
main() {
int a[1];
    a[0] = 4;
}
$ cat hello4.c
main() {
int a[1]; /* a /* nested comment */. */
    a[0] = 4; /* an /* incorrectly nested comment in
                /* two lines */ *****/
}
$ a.out < hello4.c
main() {
int a[1];
Error
    a[0] = 4;
```

2.17. Uso de Dos Analizadores

La opción `-Pprefix` de flex cambia el prefijo por defecto `yy` para todas las variables globales y funciones. Por ejemplo `-Pfoo` cambia el nombre de `yytext` a `footext`. También cambia el nombre del fichero de salida de `lex.yy.c` a `lex.foo.c`. Sigue la lista de identificadores afectados:

```
yy_create_buffer
yy_delete_buffer
yy_flex_debug
yy_init_buffer
yy_flush_buffer
yy_load_buffer_state
yy_switch_to_buffer
yyin
yyleng
yylex
yylineno
yyout
yyrestart
yytext
yywrap
```

Desde dentro del analizador léxico puedes referirte a las variables globales y funciones por cualquiera de los nombres, pero externamente tienen el nombre cambiado. Esta opción nos permite enlazar diferentes programas flex en un mismo ejecutable.

Sigue un ejemplo de uso de dos analizadores léxicos dentro del mismo programa:

```
$ cat one.l
%%
one {printf("1\n"); return 1;}
. {printf("First analyzer: %s\n",yytext);}
%%

int onewrap(void) {
    return 1;
}

$ cat two.l
%%
two {printf("2\n"); return 2;}
. {printf("Second analyzer: %s\n",yytext);}
%%

int twowrap(void) {
    return 1;
}

$ cat onetwo.c
main() {
    onelex();
    twolex();
}
```

Como hemos mencionado, la compilación *flex* se debe realizar con el opción `-P`, que cambia el prefijo por defecto `yy` de las funciones y variables accesibles por el usuario. El mismo efecto puede conseguirse utilizando la opción `prefix`, escribiendo `%option prefix="one"` y `%option prefix="two"` en los respectivos programas `one.l` y `two.l`.

```

$ flex -Pone one.l
$ flex -Ptwo two.l
$ ls -ltr | tail -2
-rw-rw----  1 pl      casiano    36537 Nov  7 09:52 lex.one.c
-rw-rw----  1 pl      casiano    36524 Nov  7 09:52 lex.two.c
$ gcc onetwo.c lex.one.c lex.two.c
$ a.out
two
First analyzer: t
First analyzer: w
First analyzer: o

one
1
one
Second analyzer: o
Second analyzer: n
Second analyzer: e

two
2
$

```

2.18. La Opción outfile

Es posible utilizar la opción `-ooutput.c` para escribir el analizador léxico en el fichero `output.c` en vez de en `lex.yy.c`. El mismo efecto puede obtenerse usando la opción `outfile="output.c"` dentro del programa `lex`.

2.19. Leer desde una Cadena: YY_INPUT

En general, la rutina que hace el análisis léxico, `yylex()`, lee su entrada a través de la macro `YY_INPUT`. Esta macro es llamada con tres parámetros

```
YY_INPUT(buf,result,max)
```

el primero, *buf* es utilizado para guardar la entrada. el tercero *max* indica el número de caracteres que `yylex()` pretende leer de la entrada. El segundo *result* contendrá el número de caracteres realmente leídos. Para poder leer desde una cadena (*string*) basta con modificar `YY_INPUT` para que copie los datos de la cadena en el *buffer* pasado como parámetro a `YY_INPUT`. Sigue un ejemplo:

```

$ cat string.l
%{
#undef YY_INPUT
#define YY_INPUT(b,r,m) (r = yystringinput(b,m))
#define min(a,b) ((a<b)?(a):(b))
%}

%%
[0-9]+ printf("Num-");
[a-zA-Z][a-zA-Z_0-9]* printf("Id-");
[ \t]+
. printf("%c-",yytext[0]);
%%

```



```

extern char string[];
extern char *yyinputptr;
extern char *yyinputlim;

int yystringinput(char *buf, int maxsize) {
    int n = min(maxsize, yyinputlim-yyinputptr);

    if (n > 0) {
        memcpy(buf, yyinputptr, n);
        yyinputptr += n;
    }
    return n;
}

int yywrap() { return 1; }

```

Este es el fichero conteniendo la función *main*:

```

$ cat stringmain.c
char string[] = "one=1;two=2";
char *yyinputptr;
char *yyinputlim;

main() {
    yyinputptr = string;
    yyinputlim = string + strlen(string);
    yylex();
    printf("\n");
}

```

Y esta es la salida:

```

$ a.out
Id==Num-;-Id==Num-

```

La cadena `string = "one=1;two=2"` definida en la línea 2 ha sido utilizada como entrada para el análisis léxico.

2.20. El operador de “trailing context” o “lookahead” positivo

En el lenguaje FORTRAN original los “blancos” no eran significativos y no se distinguía entre mayúsculas y minúsculas. Así pues la cadena `do i = 1, 10` es equivalente a la cadena `DOI=1,10`. Un conocido conflicto ocurre entre una cadena con la estructura `do i = 1.10` (esto es `DOI=1.10`) y la cadena anterior. En la primera `DO` e `I` son dos “tokens” diferentes, el primero correspondiendo a la palabra reservada que indica un bucle. En la segunda, `DOI` constituye un único “token” y la sentencia se refiere a una asignación. El conflicto puede resolverse utilizando el operador de “trailing” `r/s`. Como se mencionó, el operador de “trailing” `r/s` permite reconocer una `r` pero sólo si va seguida de una `s`. El texto casado con `s` se incluye a la hora de decidir cual es el emparejamiento mas largo, pero se devuelve a la entrada cuando se ejecuta la acción. La acción sólo ve el texto asociado con `r`. El fichero `fortran4.1` ilustra una posible solución:

```

cat fortran4.1
%array
%{
#include <string.h>

```

```

#undef YY_INPUT
#define YY_INPUT(buf,result,max) (result = my_input(buf,max))
%}
number [0-9]+
integer [+-]?{number}
float ({integer}\.{number}?|\.{number})(E{integer})?
label [A-Z0-9]+
id [A-Z]{label}*
%%
DO/{label}={number}\, { printf("do loop\n"); }
{id} { printf("Identifier %s\n",yytext); }
{number} { printf("Num %d\n",atoi(yytext)); }
{float} { printf("Float %f\n",atof(yytext)); }
(.|\n)
%%

int my_input(char *buf, int max)
{
    char *q1, *q2, *p = (char *) malloc(max);
    int i;
    if ('\0' != fgets(p,max,yyin)) {
        for(i=0, q1=buf, q2=p;(*q2 != '\0');q2++) {
            if (*q2 != ' ') { *q1++ = toupper(*q2); i++; };
        }
        free(p);
        return i;
    }
    else exit(1);
}

```

La función

```
char *fgets(char *s, int size, FILE *stream)
```

lee a lo mas uno menos que `size` caracteres desde `stream` y los almacena en el *buffer* apuntado por `s`. La lectura termina después de un EOF o un retorno de carro. Si se lee un `\n`, se almacena en el *buffer*. La función pone un carácter nulo `\0` como último carácter en el *buffer*.

A continuación, puedes ver los detalles de una ejecución:

```

$ flex fortran4.1; gcc lex.yy.c -lfl; a.out
do j = 1 . 10
Identifier DOJ
Float 1.100000
do k = 1, 5
do loop
Identifier K
Num 1
Num 5

```

2.21. Manejo de directivas include

El análisis léxico de algunos lenguajes requiere que, durante la ejecución, se realice la lectura desde diferentes ficheros de entrada. El ejemplo típico es el manejo de las directivas *include file* existentes en la mayoría de los lenguajes de programación.

¿Donde está el problema? La dificultad reside en que los analizadores generados por *flex* proveen almacenamiento intermedio (*buffers*) para aumentar el rendimiento. No basta con reescribir nuestro

propio *YY_INPUT* de manera que tenga en cuenta con que fichero se esta trabajando. El analizador sólo llama a *YY_INPUT* cuando alcanza el final de su *buffer*, lo cual puede ocurrir bastante después de haber encontrado la sentencia *include* que requiere el cambio de fichero de entrada.

```
$ cat include.l
%x incl
%{
#define yywrap() 1
#define MAX_INCLUDE_DEPTH 10
YY_BUFFER_STATE include_stack[MAX_INCLUDE_DEPTH];
int include_stack_ptr = 0;
%}
%%
include      BEGIN(incl);
.            ECHO;
<incl>[ \t]*
<incl>[^ \t\n]+ { /* got the file name */
    if (include_stack_ptr >= MAX_INCLUDE_DEPTH) {
        fprintf(stderr,"Includes nested too deeply\n");
        exit(1);
    }
    include_stack[include_stack_ptr++] = YY_CURRENT_BUFFER;
    yyin = fopen(yytext,"r");
    if (!yyin) {
        fprintf(stderr,"File %s not found\n",yytext);
        exit(1);
    }
    yy_switch_to_buffer(yy_create_buffer(yyin, YY_BUF_SIZE));
    BEGIN(INITIAL);
}
<<EOF>> {
    if ( --include_stack_ptr < 0) {
        yyterminate();
    } else {
        yy_delete_buffer(YY_CURRENT_BUFFER);
        yy_switch_to_buffer(include_stack[include_stack_ptr]);
    }
}
%%
main(int argc, char ** argv) {

    yyin = fopen(argv[1],"r");
    yylex();
}
```

La función *yy_create_buffer*(*yyin*, *YY_BUF_SIZE*)); crea un *buffer* lo suficientemente grande para mantener *YY_BUF_SIZE* caracteres. Devuelve un *YY_BUFFER_STATE*, que puede ser pasado a otras rutinas. *YY_BUFFER_STATE* es un puntero a una estructura de datos opaca (*struct yy_buffer_state*) que contiene la información para la manipulación del *buffer*. Es posible por tanto inicializar un puntero *YY_BUFFER_STATE* usando la expresión ((*YY_BUFFER_STATE*) 0).

La función *yy_switch_to_buffer*(*YY_BUFFER_STATE new_buffer*); conmuta la entrada del analizador léxico. La función *void yy_delete_buffer*(*YY_BUFFER_STATE buffer*) se usa para recuperar la memoria consumida por un *buffer*. También se pueden limpiar los contenidos actuales de un *buffer* llamando a: *void yy_flush_buffer*(*YY_BUFFER_STATE buffer*)

La regla especial <<EOF>> indica la acción a ejecutar cuando se ha encontrado un final de fichero e `yywrap()` retorna un valor distinto de cero. Cualquiera que sea la acción asociada, esta debe terminar con uno de estos cuatro supuestos:

1. Asignar `yyin` a un nuevo fichero de entrada.
2. Ejecutar `return`.
3. Ejecutar `yyterminate()`.
4. Cambiar a un nuevo buffer usando `yy_switch_to_buffer()`.

La regla <<EOF>> no se puede mezclar con otros patrones.

Este es el resultado de una ejecución del programa:

```
$ cat hello.c
#include hello2.c
main() <%
int a<:1:>; /* a comment */
    a<:0:> = 4; /* a comment in
                two lines */
    printf("\thell\157\nworld! a(0) is %d\n",a<:0:>);
%>
$ cat hello2.c
#include hello3.c
/* file hello2.c */
$ cat hello3.c
/*
third file
*/
$ flex include.1 ; gcc lex.yy.c ; a.out hello.c
##/*
third file
*/

/* file hello2.c */

main() <%
int a<:1:>; /* a comment */
    a<:0:> = 4; /* a comment in
                two lines */
    printf("\thell\157\nworld! a(0) is %d\n",a<:0:>);
%>
```

Una alternativa a usar el patrón <<EOF>> es dejar la responsabilidad de recuperar el *buffer* anterior a `yywrap()`. En tal caso suprimiríamos esta parajea patrón-acción y reescribiríamos `yywrap()`:

```
%x incl
%{
#define MAX_INCLUDE_DEPTH 10
YY_BUFFER_STATE include_stack[MAX_INCLUDE_DEPTH];
int include_stack_ptr = 0;
%}
%%
include      BEGIN(incl);
.            ECHO;
```

```

<incl>[ \t]*
<incl>[^ \t\n]+ { /* got the file name */
    if (include_stack_ptr >= MAX_INCLUDE_DEPTH) {
        fprintf(stderr, "Includes nested too deeply\n");
        exit(1);
    }
    include_stack[include_stack_ptr++] = YY_CURRENT_BUFFER;
    yyin = fopen(yytext, "r");
    if (!yyin) {
        fprintf(stderr, "File %s not found\n", yytext);
        exit(1);
    }
    yy_switch_to_buffer(yy_create_buffer(yyin, YY_BUF_SIZE));
    BEGIN(INITIAL);
}
%%
main(int argc, char ** argv) {

    yyin = fopen(argv[1], "r");
    yylex();
}

int yywrap() {
    if ( --include_stack_ptr < 0) {
        return 1;
    } else {
        yy_delete_buffer(YY_CURRENT_BUFFER);
        yy_switch_to_buffer(include_stack[include_stack_ptr]);
        return 0;
    }
}

```

2.22. Análisis Léxico desde una Cadena: yy_scan_string

El objetivo de este ejercicio es mostrar como realizar un análisis léxico de los argumentos pasados en la línea de comandos. Para ello *flex* provee la función `yy_scan_string(const char * str)`. Esta rutina crea un nuevo *buffer* de entrada y devuelve el correspondiente manejador `YY_BUFFER_STATE` asociado con la cadena `str`. Esta cadena debe estar terminada por un carácter `\0`. Podemos liberar la memoria asociada con dicho *buffer* utilizando `yy_delete_buffer(BUFFER)`. La siguiente llamada a `yylex()` realizará el análisis léxico de la cadena `str`.

```

$ cat scan_str.l
%%
[0-9]+    printf("num\n");
[a-zA-Z]+ printf("Id\n");
%%
main(int argc, char ** argv) {
    int i;

    for(i=1; i<argc; i++) {
        yy_scan_string(argv[i]);
        yylex();
        yy_delete_buffer(YY_CURRENT_BUFFER);
    }
}

```

```

    }
}

int yywrap() { return 1; }
$ flex scan_str.l ; gcc lex.yy.c ; a.out Hello World! 1234
Id
Id
!num

```

Alternativamente, la función `main()` podría haber sido escrita así:

```

main(int argc, char ** argv) {
    int i;
    YY_BUFFER_STATE p;

    for(i=1;i<argc;i++) {
        p = yy_scan_string(argv[i]);
        yylex();
        yy_delete_buffer(p);
    }
}

```

La función `yy_scan_bytes(const char * bytes, int len)` hace lo mismo que `yy_scan_string` pero en vez de una cadena terminada en el carácter nulo, se usa la longitud `len`. Ambas funciones `yy_scan_string(const char * str)` y `yy_scan_bytes(const char * bytes, int len)` hacen una copia de la cadena pasada como argumento.

Estas dos funciones crean una copia de la cadena original. Es mejor que sea así, ya que `yylex()` modifica los contenidos del *buffer* de trabajo. Si queremos evitar la copia, podemos usar

```
yy_scan_buffer(char *base, yy_size_t size),
```

la cual trabaja directamente con el *buffer* que comienza en `base`, de tamaño `size bytes`, los últimos dos de los cuáles deben ser `YY_END_OF_BUFFER_CHAR` (ASCII NUL). Estos dos últimos *bytes* no son “escaneados”. El área de rastreo va desde `base[0]` a `base[size-2]`, inclusive. Si nos olvidamos de hacerlo de este modo y no establecemos los dos *bytes* finales, la función `yy_scan_buffer()` devuelve un puntero nulo y no llega a crear el nuevo buffer de entrada. El tipo `yy_size_t` es un tipo entero. Como cabe esperar, `size` se refiere al tamaño del *buffer*.

2.23. Análisis de la Línea de Comandos y 2 Analizadores

El objetivo de este ejercicio es mostrar como realizar un análisis léxico de los argumentos pasados en la línea de comandos. Para ello diseñaremos una librería que proporcionará un función `yylexarg(argc, argv)` que hace el análisis de la línea de acuerdo con la especificación *flex* correspondiente. En el ejemplo, esta descripción del analizador léxico es proporcionada en el fichero *fl.l*. Para complicar un poco mas las cosas, supondremos que queremos hacer el análisis léxico de un fichero (especificado en la línea de comandos) según se especifica en un segundo analizador léxico *trivial.l*. El siguiente ejemplo de ejecución muestra la conducta del programa:

```

$ fl -v -V -f tokens.h
verbose mode is on
version 1.0
File name is: tokens.h
Analyzing tokens.h
#-id-blanks-id-blanks-int-blanks-#-id-blanks-id-blanks-int-blanks-#-id-blanks-id-blanks
-int-blanks-#-id-blanks-id-blanks-int-blanks-#-id-blanks-id-blanks-int-blanks-

```

Los contenidos del fichero *Makefile* definen las dependencias y la estructura de la aplicación:

```

$ cat Makefile
LIBS=-lflarg
CC=gcc -g
LIBPATH=-L. -L~/lib
INCLUDES=-I. -I~/include

fl: main.c lex.arg.c lex.yy.c libflarg.a tokens.h
    $(CC) $(LIBPATH) $(INCLUDES) main.c lex.arg.c lex.yy.c $(LIBS) -o fl
lex.arg.c: fl.l
    flex -Parg fl.l
lex.yy.c: trivial.l tokens.h
    flex trivial.l
libflarg.a: flarg.o
    ar r libflarg.a flarg.o
flarg.o: flarg.c
    $(CC) -c flarg.c
clean:
$ make clean;make
rm lex.arg.c lex.yy.c *.o fl
flex -Parg fl.l
flex trivial.l
gcc -g -c flarg.c
ar r libflarg.a flarg.o
gcc -g -L. -L~/lib -I. -I~/include main.c lex.arg.c lex.yy.c -lflarg -o fl

```

Observa el uso de la opción `-Parg` en la traducción del fichero *fl.l*. Así no solo el fichero generado por *flex*, sino todas las variables y rutinas accesibles estarán prefijadas por *arg* en vez de *yy*. La librería la denominamos *libflarg.a*. (*flarg* por flex arguments). El correspondiente fichero cabecera será *flarg.h*. Los fuentes de las rutinas que compondrán la librería se mantendrán en el fichero *flarg.c*.

Lo que haremos será redefinir `YY_INPUT(buf, result, max)` para que lea su entrada desde la línea de argumentos.

```

$ cat flarg.h
int yyarglex(int argc, char **argv);
int YY_input_from_argv(char *buf, int max);
int argwrap(void);

#undef YY_INPUT
#define YY_INPUT(buf,result,max) (result = YY_input_from_argv(buf,max))

```

La función `int YY_input_from_argv(char *buf, int max)` utiliza los punteros `char **YY_targv` y `char **YY_arglim` para moverse a través de la familia de argumentos. Mientras que el primero es utilizado para el recorrido, el segundo marca el límite final. Su inicialización ocurre en

```
yyarglex(int argc, char **argv)
```

con las asignaciones:

```
YY_targv = argv+1;
YY_arglim = argv+argc;
```

despues, de lo cual, se llama al analizador léxico generado, *arglex* .

```

$ cat flarg.c
char **YY_targv;
char **YY_arglim;

```

```

int YY_input_from_argv(char *buf, int max)
{
    static unsigned offset = 0;

    int len, copylen;

    if (YY_targv >= YY_arglim) return 0;      /* EOF */
    len = strlen(*YY_targv)-offset;           /* amount of current arg */
    if(len >= max) {copylen = max-1; offset += copylen; }
    else copylen = len;
    if(len > 0) memcpy(buf, YY_targv[0]+offset, copylen);
    if(YY_targv[0][offset+copylen] == '\0') { /* end of arg */
        buf[copylen] = ' '; copylen++; offset = 0; YY_targv++;
    }
    return copylen;
}

int yyarglex(int argc, char **argv) {
    YY_targv = argv+1;
    YY_arglim = argv+argc;
    return arglex();
}

int argwrap(void) {
    return 1;
}

```

El fichero *fl.l* contiene el analizador léxico de la línea de comandos:

```

$ cat fl.l
%{
unsigned verbose;
unsigned thereisfile;
char *progName;
char fileName[256];
#include "flarg.h"
#include "tokens.h"
%}

%%
-h      |
"-?"    |
-help   { printf("usage is: %s [-help | -h | -? ] [-verbose | -v]"
                " [-Version | -V]"
                " [-f filename]\n", progName);
        }

-v      |
-verbose { printf("verbose mode is on\n"); verbose = 1; }

-V      |
-version { printf("version 1.0\n"); }

```



```
-f[[:blank:]]+[^ \t\n]+ {
    strcpy(fileName, argtext+3);
    printf("File name is: %s\n", fileName);
    thereisfile = 1;
}

.
```

\n

Observe el uso de la clase `[[:blank:]]` para reconocer los blancos. Las clases son las mismas que las introducidas en *gawk*.

El análisis léxico del fichero que se lee después de procesar la línea de comandos es descrito en *trivial.l*. Partiendo de *trivial.l*, la ejecución del *Makefile* da lugar a la construcción por parte de *flex* del fichero *lex.yy.c* conteniendo la rutina *yylex*.

```
$ cat trivial.l
%{
#include "tokens.h"
%}
digit [0-9]
id [a-zA-Z][a-zA-Z0-9]+
blanks [ \t\n]+
operator [+*/*-]
%%
{digit}+ {return INTTOKEN; }
{digit}+"."{digit}+ {return FLOATTOKEN; }
{id} {return IDTOKEN;}
{operator} {return OPERATORTOKEN;}
{blanks} {return BLANKTOKEN;}
. {return (int) yytext[0];}
%%
int yywrap() {
    return 1;
}
```

El fichero *tokens.h* contiene la definición de los *tokens* y es compartido con *main.c*.

```
$ cat tokens.h
#define INTTOKEN 256
#define FLOATTOKEN 257
#define IDTOKEN 258
#define OPERATORTOKEN 259
#define BLANKTOKEN 260
```

Nos queda por presentar el fichero *main.c*:

```
$ cat main.c
#include <stdio.h>
#include "flarg.h"
#include "tokens.h"
extern unsigned verbose;
extern unsigned thereisfile;
extern char *progName;
extern char fileName[256];
extern FILE *yyin;
```

```

main(int argc, char **argv) {
    unsigned lookahead;
    FILE * file;

    progName = *argv;
    yyarglex(argc,argv);
    if (thereisfile) {
        if (verbose) printf("Analyzing %s\n",fileName);
        file = (fopen(fileName,"r"));
        if (file == NULL) exit(1);
        yyin = file;
        while (lookahead = yylex()) {
            switch (lookahead) {
                case INTTOKEN:
                    printf("int-");
                    break;
                case FLOATTOKEN:
                    printf("float-");
                    break;
                case IDTOKEN:
                    printf("id-");
                    break;
                case OPERATORTOKEN:
                    printf("operator-");
                    break;
                case BLANKTOKEN:
                    printf("blanks-");
                    break;
                default: printf("%c-",lookahead);
            }
        } /* while */
        printf("\n");
    } /* if */
}

```

2.24. Declaraciones pointer y array

Como se comentó, las opciones `%pointer` y `%array` controlan la definición que *flex* hace de *yytext*. en el caso en que eligamos la opción `%array` la variable `YYLMAX` controla el tamaño del *array*. Supongamos que en el fichero *trivial.l* del ejemplo anterior introducimos las siguientes modificaciones:

```

$ cat trivial.l
%array
%{
#undef YYLMAX
#define YYLMAX 4
#include "tokens.h"
%}
digit [0-9]
id [a-zA-Z][a-zA-Z0-9]+
blanks [ \t\n]+
operator [+*/*-]
%%
{digit}+ {return INTTOKEN; }

```

```

{digit}+"."{digit}+ {return FLOATTOKEN; }
{id} {return IDTOKEN;}
{operator} {return OPERATORTOKEN;}
{blanks} {return BLANKTOKEN;}
. {return (int) yytext[0];}
%%
int yywrap() {
    return 1;
}

```

En tal caso, la definición excesivamente pequeña de YYLMAX provoca un error en tiempo de ejecución:

```

$ fl -V -f tokens.h
version 1.0
File name is: tokens.h
token too large, exceeds YYLMAX

```

2.25. Las Macros YY_USER_ACTION, yy_act e YY_NUM_RULES

La macro YY_USER_ACTION permite ejecutar una acción inmediatamente después del “emparejamiento” y antes de la ejecución de la acción asociada. cuando se la invoca, la variable yy_act contiene el número de la regla que ha emparejado (las reglas se numeran a partir de uno). La macro YY_NUM_RULES contiene el número de reglas, incluyendo la regla por defecto.

El siguiente programa aprovecha dichas macros para mostrar las frecuencias de uso de las reglas.

```

$ cat user_action.l
%array
%{
#include <string.h>

int ctrl[YY_NUM_RULES];
#undef YY_USER_ACTION
#define YY_USER_ACTION { ++ctrl[yy_act]; }
%}
number [0-9]+
id      [a-zA-Z_]+[a-zA-Z0-9_]*
whites  [ \t\n]+
%%
{id}
{number}
{whites}
.
%%

int yywrap() {
    int i;

    for(i=1;i<YY_NUM_RULES;i++)
        printf("Rule %d: %d occurrences\n",i,ctrl[i]);
}

$ flex user_action.l ; gcc lex.yy.c -lfl ; a.out
a=b+2*(c-4)
Rule 1: 3 occurrences

```

Rule 2: 2 occurrences
Rule 3: 1 occurrences
Rule 4: 6 occurrences

2.26. Las opciones interactive

La opción `option always-interactive` hace que `flex` genere un analizador que considera que su entrada es “interactiva”. Concretamente, el analizador para cada nuevo fichero de entrada, intenta determinar si se trata de un a entrada interactiva o desde fichero haciendo una llamada a la función `isatty()`. Vea un ejemplo de uso de esta función:

```
$ cat isatty.c
#include <unistd.h>
#include <stdio.h>
main() {

    if (isatty(0))
        printf("interactive\n");
    else
        printf("non interactive\n");
}
$ gcc isatty.c; a.out
interactive
$ a.out < isatty.c
non interactive
$
```

cuando se usa la opción `option always-interactive`, se elimina esta llamada.

2.27. La macro YY_BREAK

Las acciones asociadas con los patrones se agrupan en la rutina de análisis léxico `yylex()` en una sentencia `switch` y se separan mediante llamadas a la macro `YY_BREAK`. Así, al compilar con `flex` el siguiente fichero `.l`

```
$ cat interactive.l
%%
. printf("::%c",yytext[0]);
\n printf("::%c",yytext[0]);
```

tenemos el fichero de salida `lex.yy.c` que aparece a continuación (hemos omitido las líneas de código en las que estamos menos interesados, sustituyéndolas por puntos suspensivos)

```
/* A lexical scanner generated by flex */
....
#define YY_NUM_RULES 3
#line 1 "interactive.l"
#define INITIAL 0
#line 363 "lex.yy.c"
....
YY_DECL {
    ....
#line 1 "interactive.l"
#line 516 "lex.yy.c"
```

```

....
if ( yy_init ) {
    yy_init = 0;
#ifdef YY_USER_INIT
    YY_USER_INIT;
#endif
    if ( ! yy_start ) yy_start = 1; /* first start state */
    if ( ! yyin ) yyin = stdin;
    if ( ! yyout ) yyout = stdout;
    if ( ! yy_current_buffer ) yy_current_buffer = yy_create_buffer( yyin, YY_BUF_SIZE );
    yy_load_buffer_state();
}
while ( 1 ) /* loops until end-of-file is reached */ {
    .....
yy_match:
    do {
        .....
    }
    .....
yy_find_action:
    .....
    YY_DO_BEFORE_ACTION;
do_action: /* This label is used only to access EOF actions. */
    switch ( yy_act ) { /* beginning of action switch */
        case 0:
            .....
            goto yy_find_action;
        case 1:
            YY_RULE_SETUP
            #line 2 "interactive.l"
            printf(":::%c",yytext[0]);
            YY_BREAK
        case 2:
            YY_RULE_SETUP
            #line 3 "interactive.l"
            printf(":::%c",yytext[0]);
            YY_BREAK
        case 3:
            YY_RULE_SETUP
            #line 4 "interactive.l"
            ECHO;
            YY_BREAK
        #line 614 "lex.yy.c"
        case YY_STATE_EOF(INITIAL):
            yyterminate();
        case YY_END_OF_BUFFER:
            { ..... }
        default:
            YY_FATAL_ERROR("fatal flex scanner internal error--no action found");
    } /* end of action switch */
} /* end of scanning one token */
} /* end of yylex */

```

```
#if YY_MAIN
int main()
{ yylex(); return 0; }
#endif
#line 4 "interactive.l"
```

Por defecto, la macro `YY_BREAK` es simplemente un `break`. Si cada acción de usuario termina en un `return`, puedes encontrarte con que el compilador genera un buen número de `warning! unreachable code`. Puedes entonces redefinir `YY_BREAK` a vacío y evitar estos mensajes.

Capítulo 3

Expresiones Regulares en Perl

3.1. Introducción

Los rudimentos de las expresiones regulares pueden encontrarse en los trabajos pioneros de McCulloch y Pitts (1940) sobre redes neuronales. El lógico Stephen Kleene definió formalmente el álgebra que denominó *conjuntos regulares* y desarrolló una notación para la descripción de dichos conjuntos, las *expresiones regulares*.

Durante las décadas de 1960 y 1970 hubo un desarrollo formal de las expresiones regulares. Una de las primeras publicaciones que utilizan las expresiones regulares en un marco informático es el artículo de 1968 de Ken Thompson *Regular Expression Search Algorithm* en el que describe un compilador de expresiones regulares que produce código objeto para un IBM 7094. Este compilador dio lugar al editor *qed*, en el cual se basó el editor de Unix *ed*. Aunque las expresiones regulares de este último no eran tan sofisticadas como las de *qed*, fueron las primeras en ser utilizadas en un contexto no académico. Se dice que el comando global **g** en su formato **g/re/p** que utilizaba para imprimir (opción **p**) las líneas que casan con la expresión regular **re** dio lugar a un programa separado al que se denominó *grep*.

Las expresiones regulares facilitadas por las primeras versiones de estas herramientas eran limitadas. Por ejemplo, se disponía del cierre de Kleene ***** pero no del cierre positivo **+** o del operador opcional **?**. Por eso, posteriormente, se han introducido los metacaracteres **\+** y **\?**. Existían numerosas limitaciones en dichas versiones, por ej. **\$** sólo significa “final de línea” al final de la expresión regular. Eso dificulta expresiones como

```
grep 'cierre$\|^Las' viq.tex
```

Sin embargo, la mayor parte de las versiones actuales resuelven correctamente estos problemas:

```
nereida:~/viq> grep 'cierre$\|^Las' viq.tex
```

Las expresiones regulares facilitadas por las primeras versiones de estas herramientas eran limitadas. Por ejemplo, se disponía del cierre de Kleene **\verb|*|** pero no del cierre

```
nereida:~/viq>
```

De hecho AT&T Bell labs añadió numerosas funcionalidades, como por ejemplo, el uso de **\{min, max\}**, tomada de *lex*. Por esa época, Alfred Aho escribió *egrep* que, no sólo proporciona un conjunto más rico de operadores sino que mejoró la implementación. Mientras que el *grep* de Ken Thompson usaba un autómata finito no determinista (NFA), la versión de *egrep* de Aho usa un autómata finito determinista (DFA).

En 1986 Henry Spencer desarrolló la librería *regex* para el lenguaje C, que proporciona un conjunto consistente de funciones que permiten el manejo de expresiones regulares. Esta librería ha contribuido a “homogeneizar” la sintaxis y semántica de las diferentes herramientas que utilizan expresiones regulares (como *awk*, *lex*, *sed*, ...).

Véase También

- La sección *Expresiones Regulares en Otros Lenguajes* 3.3

- Regular Expressions Cookbook. Jan Goyvaerts, Steven Levithan
- PCRE (Perl Compatible Regular Expressions) en la Wikipedia
- PCRE (Perl Compatible Regular Expressions)
- Java Regular Expressions
- C# Regular Expressions
- .NET Framework Regular Expressions

3.1.1. Un ejemplo sencillo

Matching en Contexto Escalar

```
pl@nereida:~/Lperltesting$ cat -n c2f.pl
1  #!/usr/bin/perl -w
2  use strict;
3
4  print "Enter a temperature (i.e. 32F, 100C):\n";
5  my $input = <STDIN>;
6  chomp($input);
7
8  if ($input !~ m/^[(-+)?[0-9]+(\.[0-9]*)?)\s*([CF])$/i) {
9      warn "Expecting a temperature, so don't understand \"$input\".\n";
10 }
11 else {
12     my $InputNum = $1;
13     my $type = $3;
14     my ($celsius, $fahrenheit);
15     if ($type eq "C" or $type eq "c") {
16         $celsius = $InputNum;
17         $fahrenheit = ($celsius * 9/5)+32;
18     }
19     else {
20         $fahrenheit = $InputNum;
21         $celsius = ($fahrenheit -32)*5/9;
22     }
23     printf "%.2f C = %.2f F\n", $celsius, $fahrenheit;
24 }
```

Véase también:

- `perldoc perlrequick`
- `perldoc perlretut`
- `perldoc perlre`
- `perldoc perlref`

Ejecución con el depurador:

```
pl@nereida:~/Lperltesting$ perl -wd c2f.pl
Loading DB routines from perl5db.pl version 1.28
Editor support available.
```



```

Enter h or 'h h' for help, or 'man perldebug' for more help.
main::(c2f.pl:4):      print "Enter a temperature (i.e. 32F, 100C):\n";
DB<1> c 8
Enter a temperature (i.e. 32F, 100C):
32F
main::(c2f.pl:8):      if ($input !~ m/^( [-+]?[0-9]+(\.[0-9]*)?)\s*([CF])$/i) {
DB<2> n
main::(c2f.pl:12):     my $InputNum = $1;
DB<2> x ($1, $2, $3)
0 32
1 undef
2 'F'
DB<3> use YAPE::Regex::Explain
DB<4> p YAPE::Regex::Explain->new('([-+]?[0-9]+(\.[0-9]*)?)\s*([CF])$')->explain
The regular expression:
(?-imsx:([-+]?[0-9]+(\.[0-9]*)?)\s*([CF])$)
matches as follows:

```

NODE	EXPLANATION
(?-imsx:	group, but do not capture (case-sensitive) (with ^ and \$ matching normally) (with . not matching \n) (matching whitespace and # normally):
(group and capture to \1:
[-+]?	any character of: '-', '+' (optional (matching the most amount possible))
[0-9]+	any character of: '0' to '9' (1 or more times (matching the most amount possible))
(group and capture to \2 (optional (matching the most amount possible)):
\.	'.'
[0-9]*	any character of: '0' to '9' (0 or more times (matching the most amount possible))
)?	end of \2 (NOTE: because you're using a quantifier on this capture, only the LAST repetition of the captured pattern will be stored in \2)
)	end of \1
\s*	whitespace (\n, \r, \t, \f, and " ") (0 or more times (matching the most amount possible))

(group and capture to \3:
[CF]	any character of: 'C', 'F'
)	end of \3
\$	before an optional \n, and the end of the string
)	end of grouping

Dentro de una expresión regular es necesario referirse a los textos que casan con el primer, paréntesis, segundo, etc. como \1, \2, etc. La notación \$1 se refiere a lo que casó con el primer paréntesis en el último *matching*, no en el actual. Veamos un ejemplo:

```
pl@nereida:~/Lperltesting$ cat -n dollar1slash1.pl
1    #!/usr/bin/perl -w
2    use strict;
3
4    my $a = "hola juanito";
5    my $b = "adios anita";
6
7    $a =~ /(ani)/;
8    $b =~ s/(adios) *($1)/\U$1 $2/;
9    print "$b\n";
```

Observe como el \$1 que aparece en la cadena de reemplazo (línea 8) se refiere a la cadena **adios** mientras que el \$1 en la primera parte contiene **ani**:

```
pl@nereida:~/Lperltesting$ ./dollar1slash1.pl
ADIOS ANIta
```

Ejercicio 3.1.1. *Indique cuál es la salida del programa anterior si se sustituye la línea 8 por*

```
$b =~ s/(adios) *(\1)/\U$1 $2/;
```

Número de Paréntesis

El número de paréntesis con memoria no está limitado:

```
pl@nereida:~/Lperltesting$ perl -wde 0
main::(-e:1): 0
123456789ABCDEF
DB<1> $x = "123456789AAAAAA"
1 2 3 4 5 6 7 8 9 10 11 12
DB<2> $r = $x =~ /(.) (.) (.) (.) (.) (.) (.) (.) (.) (.) \11/; print "$r\n$10\n$11\n"
1
A
A
```

Véase el siguiente párrafo de **perlre** (sección Capture buffers):

There is no limit to the number of captured substrings that you may use. However Perl also uses \10, \11, etc. as aliases for \010, \011, etc. (Recall that 0 means octal, so \011 is the character at number 9 in your coded character set; which would be the 10th

character, a horizontal tab under ASCII.) Perl resolves this ambiguity by interpreting \10 as a backreference only if at least 10 left parentheses have opened before it. Likewise \11 is a backreference only if at least 11 left parentheses have opened before it. And so on. \1 through \9 are always interpreted as backreferences.

Contexto de Lista

Si se utiliza en un contexto que requiere una lista, el “pattern match” retorna una lista consistente en las subexpresiones casadas mediante los paréntesis, esto es \$1, \$2, \$3, Si no hubiera emparejamiento se retorna la lista vacía. Si lo hubiera pero no hubieran paréntesis se retorna la lista (\$&).

```
pl@nereida:~/src/perl/perltesting$ cat -n escapes.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3
 4  my $foo = "one two three four five\nsix seven";
 5  my ($F1, $F2, $Etc) = ($foo =~ /\s*(\S+)\s+(\S+)\s*(.*)/);
 6  print "List Context: F1 = $F1, F2 = $F2, Etc = $Etc\n";
 7
 8  # This is 'almost' the same than:
 9  ($F1, $F2, $Etc) = split(/\s+/, $foo, 3);
10  print "Split: F1 = $F1, F2 = $F2, Etc = $Etc\n";
```

Observa el resultado de la ejecución:

```
pl@nereida:~/src/perl/perltesting$ ./escapes.pl
List Context: F1 = one, F2 = two, Etc = three four five
Split: F1 = one, F2 = two, Etc = three four five
six seven
```

El modificador s

La opción `s` usada en una regexp hace que el punto `.` case con el retorno de carro:

```
pl@nereida:~/src/perl/perltesting$ perl -wd ./escapes.pl
main: (./escapes.pl:4): my $foo = "one two three four five\nsix seven";
DB<1> c 9
List Context: F1 = one, F2 = two, Etc = three four five
main: (./escapes.pl:9): ($F1, $F2, $Etc) = split(' ', $foo, 3);
DB<2> ($F1, $F2, $Etc) = ($foo =~ /\s*(\S+)\s+(\S+)\s*(.*)/s)
DB<3> p "List Context: F1 = $F1, F2 = $F2, Etc = $Etc\n"
List Context: F1 = one, F2 = two, Etc = three four five
six seven
```

La opción `/s` hace que `.` se empareje con un `\n`. Esto es, casa con cualquier carácter.

Veamos otro ejemplo, que imprime los nombres de los ficheros que contienen cadenas que casan con un patrón dado, incluso si este aparece disperso en varias líneas:

```
1  #!/usr/bin/perl -w
2  #use:
3  #smodifier.pl 'expr' files
4  #prints the names of the files that match with the give expr
5  undef $/; # input record separator
6  my $what = shift @ARGV;
7  while(my $file = shift @ARGV) {
8      open(FILE, "<$file");
```

```

9     $line = <FILE>;
10    if ($line =~ /$what/s) {
11        print "$file\n";
12    }
13 }
```

Ejemplo de uso:

```

> smodifier.pl 'three.*three' double.in split.pl doublee.pl
double.in
doublee.pl
```

Vea la sección 3.4.2 para ver los contenidos del fichero `double.in`. En dicho fichero, el patrón `three.*three` aparece repartido entre varias líneas.

El modificador `m`

El modificador `s` se suele usar conjuntamente con el modificador `m`. He aquí lo que dice la sección *Using character classes* de la sección 'Using-character-classes' en `perlretut` al respecto:

- *m modifier (//m): Treat string as a set of multiple lines. '.' matches any character except \n. ^ and \$ are able to match at the start or end of any line within the string.*
 - *both s and m modifiers (//sm): Treat string as a single long line, but detect multiple lines. '.' matches any character, even \n. ^ and \$, however, are able to match at the start or end of any line within the string.*
- Here are examples of //s and //m in action:*

1. `$x = "There once was a girl\nWho programmed in Perl\n";`
- 2.
3. `$x =~ /^Who/; # doesn't match, "Who" not at start of string`
4. `$x =~ /^Who/s; # doesn't match, "Who" not at start of string`
5. `$x =~ /^Who/m; # matches, "Who" at start of second line`
6. `$x =~ /^Who/sm; # matches, "Who" at start of second line`
- 7.
8. `$x =~ /girl.Who/; # doesn't match, "." doesn't match "\n"`
9. `$x =~ /girl.Who/s; # matches, "." matches "\n"`
10. `$x =~ /girl.Who/m; # doesn't match, "." doesn't match "\n"`
11. `$x =~ /girl.Who/sm; # matches, "." matches "\n"`

Most of the time, the default behavior is what is wanted, but //s and //m are occasionally very useful. If //m is being used, the start of the string can still be matched with \A and the end of the string can still be matched with the anchors \Z (matches both the end and the newline before, like \$), and \z (matches only the end):

1. `$x =~ /^Who/m; # matches, "Who" at start of second line`
2. `$x =~ /\AWho/m; # doesn't match, "Who" is not at start of string`
- 3.
4. `$x =~ /girl$/m; # matches, "girl" at end of first line`
5. `$x =~ /girl\Z/m; # doesn't match, "girl" is not at end of string`
- 6.
7. `$x =~ /Perl\Z/m; # matches, "Perl" is at newline before end`
8. `$x =~ /Perl\z/m; # doesn't match, "Perl" is not at end of string`

Normalmente el carácter `^` casa solamente con el comienzo de la cadena y el carácter `$` con el final. Los `\n` empotrados no casan con `^` ni `$`. El modificador `/m` modifica esta conducta. De este modo `^` y `$` casan con cualquier frontera de línea interna. Las anclas `\A` y `\Z` se utilizan entonces para casar con el comienzo y final de la cadena. Véase un ejemplo:

```

nereida:~/perl/src> perl -de 0
DB<1> $a = "hola\npedro"
DB<2> p "$a"
hola
pedro
DB<3> $a =~ s/./x/m
DB<4> p $a
x
pedro
DB<5> $a =~ s/^pedro$/juan/
DB<6> p "$a"
x
pedro
DB<7> $a =~ s/^pedro$/juan/m
DB<8> p "$a"
x
juan

```

El conversor de temperaturas reescrito usando contexto de lista

Reescribamos el ejemplo anterior usando un contexto de lista:

```

casiano@milllo:~/Lperltesting$ cat -n c2f_list.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3
 4  print "Enter a temperature (i.e. 32F, 100C):\n";
 5  my $input = <STDIN>;
 6  chomp($input);
 7
 8  my ($InputNum, $type);
 9
10  ($InputNum, $type) = $input =~ m/^(
11                                ([-+]?[0-9]+(?:\.[0-9]*)?) # Temperature
12                                \s*
13                                ([cCfF]) # Celsius or Farenheit
14                                $/x;
15
16  die "Expecting a temperature, so don't understand \"$input\".\n" unless defined($InputNum);
17
18  my ($celsius, $fahrenheit);
19  if ($type eq "C" or $type eq "c") {
20      $celsius = $InputNum;
21      $fahrenheit = ($celsius * 9/5)+32;
22  }
23  else {
24      $fahrenheit = $InputNum;
25      $celsius = ($fahrenheit -32)*5/9;
26  }
27  printf "%.2f C = %.2f F\n", $celsius, $fahrenheit;

```

La opción x

La opción /x en una regexp permite utilizar comentarios y espacios dentro de la expresión regular. Los espacios dentro de la expresión regular dejan de ser significativos. Si quieres conseguir un espacio

que sea significativo, usa `\s` o bien escápalo. Véase la sección 'Modifiers' en `perlre` y la sección 'Building-a-regexp' en `perlretut`.

Paréntesis sin memoria

La notación `(?: ...)` se usa para introducir paréntesis de agrupamiento sin memoria. `(?: ...)` Permite agrupar las expresiones tal y como lo hacen los paréntesis ordinarios. La diferencia es que no "memorizan" esto es no guardan nada en `$1`, `$2`, etc. Se logra así una compilación mas eficiente. Veamos un ejemplo:

```
> cat groupingpar.pl
#!/usr/bin/perl

my $a = shift;

$a =~ m/(?:hola )*(juan)/;
print "$1\n";
nereida:~/perl/src> groupingpar.pl 'hola juan'
juan
```

Interpolación en los patrones: La opción `/o`

El patrón regular puede contener variables, que serán interpoladas (en tal caso, el patrón será re-compilado). Si quieres que dicho patrón se compile una sólo vez, usa la opción `/o`.

```
pl@nereida:~/Lperltesting$ cat -n mygrep.pl
1  #!/usr/bin/perl -w
2  my $what = shift @ARGV || die "Usage $0 regexp files ...\n";
3  while (<>) {
4      print "File $ARGV, rel. line $.: $_" if (/ $what /o); # compile only once
5  }
6
```

Sigue un ejemplo de ejecución:

```
pl@nereida:~/Lperltesting$ ./mygrep.pl
Usage ./mygrep.pl regexp files ...
pl@nereida:~/Lperltesting$ ./mygrep.pl if labels.c
File labels.c, rel. line 7:          if (a < 10) goto LABEL;
```

El siguiente texto es de la sección 'Using-regular-expressions-in-Perl' en `perlretut`:

If \$pattern won't be changing over the lifetime of the script, we can add the `/o` modifier, which directs Perl to only perform variable substitutions once

Otra posibilidad es hacer una compilación previa usando el operador `qr` (véase la sección 'Regexp-Quote-Like-Operators' en `perlop`). La siguiente variante del programa anterior también compila el patrón una sólo vez:

```
pl@nereida:~/Lperltesting$ cat -n mygrep2.pl
1  #!/usr/bin/perl -w
2  my $what = shift @ARGV || die "Usage $0 regexp files ...\n";
3  $what = qr{$what};
4  while (<>) {
5      print "File $ARGV, rel. line $.: $_" if (/ $what /);
6  }
```

Véase

- El nodo en `perlmonks` `/o is dead, long live qr//!` por diotalevi

Cuantificadores greedy

El siguiente extracto de la sección *Matching Repetitions* en la sección 'Matching-repetitions' en `perlretut` ilustra la semántica *greedy* de los operadores de repetición `*+{}?` etc.

For all of these quantifiers, Perl will try to match as much of the string as possible, while still allowing the regexp to succeed. Thus with `/a?...`, Perl will first try to match the regexp with the `a` present; if that fails, Perl will try to match the regexp without the `a` present. For the quantifier ``, we get the following:*

```
1. $x = "the cat in the hat";
2. $x =~ /^(.*)(cat)(.*)$/; # matches,
3. # $1 = 'the '
4. # $2 = 'cat'
5. # $3 = ' in the hat'
```

Which is what we might expect, the match finds the only cat in the string and locks onto it. Consider, however, this regexp:

```
1. $x =~ /^(.)(at)(.*)$/; # matches,
2. # $1 = 'the cat in the h'
3. # $2 = 'at'
4. # $3 = '' (0 characters match)
```

One might initially guess that Perl would find the `at` in `cat` and stop there, but that wouldn't give the longest possible string to the first quantifier `.`. Instead, the first quantifier `.*` grabs as much of the string as possible while still having the regexp match. In this example, that means having the `at` sequence with the final `at` in the string.*

The other important principle illustrated here is that when there are two or more elements in a regexp, the leftmost quantifier, if there is one, gets to grab as much the string as possible, leaving the rest of the regexp to fight over scraps. Thus in our example, the first quantifier `.` grabs most of the string, while the second quantifier `.*` gets the empty string. Quantifiers that grab as much of the string as possible are called maximal match or greedy quantifiers.*

When a regexp can match a string in several different ways, we can use the principles above to predict which way the regexp will match:

- **Principle 0:** Taken as a whole, any regexp will be matched at the earliest possible position in the string.
- **Principle 1:** In an alternation `a|b|c...`, the leftmost alternative that allows a match for the whole regexp will be the one used.
- **Principle 2:** The maximal matching quantifiers `?`, `*`, `+` and `{n,m}` will in general match as much of the string as possible while still allowing the whole regexp to match.
- **Principle 3:** If there are two or more elements in a regexp, the leftmost greedy quantifier, if any, will match as much of the string as possible while still allowing the whole regexp to match. The next leftmost greedy quantifier, if any, will try to match as much of the string remaining available to it as possible, while still allowing the whole regexp to match. And so on, until all the regexp elements are satisfied.

Regexp y Bucles Infinitos

El siguiente párrafo está tomado de la sección 'Repeated-Patterns-Matching-a-Zero-length-Substring' en `perlre`:

Regular expressions provide a terse and powerful programming language. As with most other power tools, power comes together with the ability to wreak havoc.

A common abuse of this power stems from the ability to make infinite loops using regular expressions, with something as innocuous as:

```
1. 'foo' =~ m{ ( o? )* }x;
```

The `o?` matches at the beginning of `'foo'` , and since the position in the string is not moved by the match, `o?` would match again and again because of the `*` quantifier.

Another common way to create a similar cycle is with the looping modifier `//g` :

```
1. @matches = ( 'foo' =~ m{ o? }xg );
```

or

```
1. print "match: <$&>\n" while 'foo' =~ m{ o? }xg;
```

or the loop implied by `split()`.

... Perl allows such constructs, by forcefully breaking the infinite loop. The rules for this are different for lower-level loops given by the greedy quantifiers `++{ }` , and for higher-level ones like the `/g` modifier or `split()` operator.

The lower-level loops are interrupted (that is, the loop is broken) when Perl detects that a repeated expression matched a zero-length substring. Thus

```
1. m{ (?: NON_ZERO_LENGTH | ZERO_LENGTH )* }x;
```

is made equivalent to

```
1. m{ (?: NON_ZERO_LENGTH )*
2. |
3. (?: ZERO_LENGTH )?
4. }x;
```

The higher level-loops preserve an additional state between iterations: whether the last match was zero-length. To break the loop, the following match after a zero-length match is prohibited to have a length of zero. This prohibition interacts with backtracking (see *Backtracking*), and so the second best match is chosen if the best match is of zero length.

For example:

```
1. $_ = 'bar';
2. s/\w??/<$&>/g;
```

results in `<><><a><><r><>` . At each position of the string the best match given by non-greedy `??` is the zero-length match, and the second best match is what is matched by `\w` . Thus zero-length matches alternate with one-character-long matches.

Similarly, for repeated `m/()/g` the second-best match is the match at the position one notch further in the string.

The additional state of being matched with zero-length is associated with the matched string, and is reset by each assignment to `pos()` . Zero-length matches at the end of the previous match are ignored during `split` .

Ejercicio 3.1.2. ■ Explique la conducta del siguiente matching:


```
DB<25> $c = 0
```

```
DB<26> print(($c++).": <$&>\n") while 'aaaabababab' =~ /a*(ab)*/g;
0: <aaaa>
1: <>
2: <a>
3: <>
4: <a>
5: <>
6: <a>
7: <>
8: <>
```

Cuantificadores *lazy*

Las expresiones *lazy* o *no greedy* hacen que el NFA se detenga en la cadena mas corta que casa con la expresión. Se denotan como sus análogas *greedy* añadiéndole el postfijo ?:

- {n,m}?
- {n,}??
- {n}??
- *?
- +?
- ??

Repasemos lo que dice la sección Matching Repetitions en la sección 'Matching-repetitions' en perlretut:

*Sometimes greed is not good. At times, we would like quantifiers to match a minimal piece of string, rather than a maximal piece. For this purpose, Larry Wall created the minimal match or non-greedy quantifiers ??, *?, +?, and {}?. These are the usual quantifiers with a ? appended to them. They have the following meanings:*

- *a?? means: match 'a' 0 or 1 times. Try 0 first, then 1.*
- *a*? means: match 'a' 0 or more times, i.e., any number of times, but as few times as possible*
- *a+? means: match 'a' 1 or more times, i.e., at least once, but as few times as possible*
- *a{n,m}? means: match at least n times, not more than m times, as few times as possible*
- *a{n,}? means: match at least n times, but as few times as possible*
- *a{n}? means: match exactly n times. Because we match exactly n times, an? is equivalent to an and is just there for notational consistency.*

Let's look at the example above, but with minimal quantifiers:

```
1. $x = "The programming republic of Perl";
2. $x =~ /^(.+?)(e|r)(.*)$/; # matches,
3. # $1 = 'Th'
4. # $2 = 'e'
5. # $3 = ' programming republic of Perl'
```

The minimal string that will allow both the start of the string `^` and the alternation to match is `Th`, with the alternation `e|r` matching `e`. The second quantifier `.*` is free to gobble up the rest of the string.

1. `$x =~ /(m{1,2}?)(.*)$/; # matches,`
2. `# $1 = 'm'`
3. `# $2 = 'ming republic of Perl'`

The first string position that this regexp can match is at the first `m` in `programming`. At this position, the minimal `m{1,2}?` matches just one `m`. Although the second quantifier `.*?` would prefer to match no characters, it is constrained by the end-of-string anchor `$` to match the rest of the string.

1. `$x =~ /(.*?)(m{1,2}?)(.*)$/; # matches,`
2. `# $1 = 'The progra'`
3. `# $2 = 'm'`
4. `# $3 = 'ming republic of Perl'`

In this regexp, you might expect the first minimal quantifier `.*?` to match the empty string, because it is not constrained by a `^` anchor to match the beginning of the word. Principle 0 applies here, however. Because it is possible for the whole regexp to match at the start of the string, it will match at the start of the string. Thus the first quantifier has to match everything up to the first `m`. The second minimal quantifier matches just one `m` and the third quantifier matches the rest of the string.

1. `$x =~ /(.*?) (m{1,2}) (.*?)$/; # matches,`
2. `# $1 = 'a'`
3. `# $2 = 'mm'`
4. `# $3 = 'ing republic of Perl'`

Just as in the previous regexp, the first quantifier `.*?` can match earliest at position `a`, so it does. The second quantifier is greedy, so it matches `mm`, and the third matches the rest of the string.

We can modify principle 3 above to take into account non-greedy quantifiers:

- **Principle 3:** If there are two or more elements in a regexp, the leftmost greedy (non-greedy) quantifier, if any, will match as much (little) of the string as possible while still allowing the whole regexp to match. The next leftmost greedy (non-greedy) quantifier, if any, will try to match as much (little) of the string remaining available to it as possible, while still allowing the whole regexp to match. And so on, until all the regexp elements are satisfied.

Ejercicio 3.1.3. Explique cuál será el resultado de el segundo comando de matching introducido en el depurador:

```
casiano@millo:~/Lperltesting$ perl -wde 0
main::(-e:1): 0
DB<1> x ('1'x34) =~ m{^(11+)\1+$}
0 11111111111111111111
DB<2> x ('1'x34) =~ m{^(11+?)\1+$}
????????????????????????????????????????
```

Descripción detallada del proceso de matching Veamos en detalle lo que ocurre durante un matching. Repasemos lo que dice la sección Matching Repetitions en la sección 'Matching-repetitions' en perlretut:

Just like alternation, quantifiers are also susceptible to backtracking. Here is a step-by-step analysis of the example

```
1. $x = "the cat in the hat";
2. $x =~ /^(.*) (at) (.*)$/; # matches,
3. # $1 = 'the cat in the h'
4. # $2 = 'at'
5. # $3 = '' (0 matches)
```

1. Start with the first letter in the string 't'.
2. The first quantifier '.*' starts out by matching the whole string 'the cat in the hat'.
3. 'a' in the regexp element 'at' doesn't match the end of the string. Backtrack one character.
4. 'a' in the regexp element 'at' still doesn't match the last letter of the string 't', so backtrack one more character.
5. Now we can match the 'a' and the 't'.
6. Move on to the third element '.*'. Since we are at the end of the string and '.*' can match 0 times, assign it the empty string.
7. We are done!

Rendimiento

La forma en la que se escribe una regexp puede dar lugar a grandes variaciones en el rendimiento. Repasemos lo que dice la sección Matching Repetitions en la sección 'Matching-repetitions' en perlretut:

Most of the time, all this moving forward and backtracking happens quickly and searching is fast. There are some pathological regexps, however, whose execution time exponentially grows with the size of the string. A typical structure that blows up in your face is of the form

```
/(a|b+)*;/;
```

The problem is the nested indeterminate quantifiers. There are many different ways of partitioning a string of length n between the $+$ and $$: one repetition with b^+ of length n , two repetitions with the first b^+ length k and the second with length $n - k$, m repetitions whose bits add up to length n , etc.*

In fact there are an exponential number of ways to partition a string as a function of its length. A regexp may get lucky and match early in the process, but if there is no match, Perl will try every possibility before giving up. So be careful with nested $$'s, $\{n,m\}$'s, and $+$'s.*

The book Mastering Regular Expressions by Jeffrey Friedl [?] gives a wonderful discussion of this and other efficiency issues.

Eliminación de Comentarios de un Programa C

El siguiente ejemplo elimina los comentarios de un programa C.

```
casiano@millor:~/Lperltesting$ cat -n comments.pl
1    #!/usr/bin/perl -w
2    use strict;
```

```

3
4   my $programe = shift @ARGV or die "Usage:\n$0 prog.c\n";
5   open(my $PROGRAM, "<$programe") || die "can't find $programe\n";
6   my $program = '';
7   {
8       local $/ = undef;
9       $program = <$PROGRAM>;
10  }
11  $program =~ s{
12      /\* # Match the opening delimiter
13      .*? # Match a minimal number of characters
14      \*/ # Match the closing delimiter
15  }[]gsx;
16
17  print $program;

```

Veamos un ejemplo de ejecución. Supongamos el fichero de entrada:

```

> cat hello.c
#include <stdio.h>
/* first
comment
*/
main() {
    printf("hello world!\n"); /* second comment */
}

```

Entonces la ejecución con ese fichero de entrada produce la salida:

```

> comments.pl hello.c
#include <stdio.h>

main() {
    printf("hello world!\n");
}

```

Veamos la diferencia de comportamiento entre `*` y `*?` en el ejemplo anterior:

```

pl@nereida:~/src/perl/perltesting$ perl5_10_1 -wde 0
main:(-e:1): 0
DB<1> use re 'debug'; 'main() /* 1c */ { /* 2c */ return; /* 3c */ }' =~ qr{(/\*.*\*/)}; pr
Compiling REx "(/\*.*\*/)"
Final program:
1: OPEN1 (3)
3: EXACT </*> (5)
5: STAR (7)
6: REG_ANY (0)
7: EXACT <*/> (9)
9: CLOSE1 (11)
11: END (0)
anchored "/" at 0 floating "/" at 2..2147483647 (checking floating) minlen 4
Guessing start of match in sv for REx "(/\*.*\*/)" against "main() /* 1c */ { /* 2c */ return;
Found floating substr "/" at offset 13...
Found anchored substr "/" at offset 7...
Starting position does not contradict /~/m...
Guessed: match at offset 7

```

```
Matching REx "(/\\*.*\\*/)" against "/* 1c */ { /* 2c */ return; /* 3c */ }"
  7 <in() > < /* 1c */ {>      | 1:OPEN1(3)
  7 <in() > < /* 1c */ {>      | 3:EXACT < /*>(5)
  9 <() /*> < 1c */ { />      | 5:STAR(7)
                                REG_ANY can match 36 times out of 2147483647...
 41 <; /* 3c > < /* }>      | 7: EXACT < /*>(9)
 43 <; /* 3c */> < }>      | 9: CLOSE1(11)
 43 <; /* 3c */> < }>      | 11: END(0)
```

Match successful!

```
/* 1c */ { /* 2c */ return; /* 3c */
Freeing REx: "(/\\*.*\\*/)"
```

```
DB<2> use re 'debug'; 'main() /* 1c */ { /* 2c */ return; /* 3c */ }' =~ qr{(/\\*.*?\\*/)}; p
Compiling REx "(/\\*.*?\\*/)"
```

Final program:

```
1: OPEN1 (3)
3: EXACT < /*> (5)
5: MINMOD (6)
6: STAR (8)
7: REG_ANY (0)
8: EXACT < /*> (10)
10: CLOSE1 (12)
12: END (0)
```

anchored "/*" at 0 floating "/*" at 2..2147483647 (checking floating) minlen 4

Guessing start of match in sv for REx "(/*.*?*/)" against "main() /* 1c */ { /* 2c */ return

Found floating substr "/*" at offset 13...

Found anchored substr "/*" at offset 7...

Starting position does not contradict /~/m...

Guessed: match at offset 7

```
Matching REx "(/\\*.*?\\*/)" against "/* 1c */ { /* 2c */ return; /* 3c */ }"
```

```
  7 <in() > < /* 1c */ {>      | 1:OPEN1(3)
  7 <in() > < /* 1c */ {>      | 3:EXACT < /*>(5)
  9 <() /*> < 1c */ { />      | 5:MINMOD(6)
  9 <() /*> < 1c */ { />      | 6:STAR(8)
                                REG_ANY can match 4 times out of 4...
 13 <* 1c > < /* { /* 2c>      | 8: EXACT < /*>(10)
 15 <1c */> < { /* 2c *>      | 10: CLOSE1(12)
 15 <1c */> < { /* 2c *>      | 12: END(0)
```

Match successful!

```
/* 1c */
```

```
Freeing REx: "(/\\*.*?\\*/)"
```

```
DB<3>
```

Véase también la documentación en la sección 'Matching-repetitions' en [perlretut](#) y la sección 'Quantifiers' en [perlre](#).

Negaciones y operadores *lazy* A menudo las expresiones $X[^X]*X$ y $X.*?X$, donde X es un carácter arbitrario se usan de forma casi equivalente.

- La primera significa:

Una cadena que no contiene X en su interior y que está delimitada por Xs

- La segunda significa:

Una cadena que comienza en X y termina en la X mas próxima a la X de comienzo

Esta equivalencia se rompe si no se cumplen las hipótesis establecidas.

En el siguiente ejemplo se intentan detectar las cadenas entre comillas dobles que terminan en el signo de exclamación:

```
pl@nereida:~/Lperltesting$ cat -n negynogreedy.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3
 4  my $b = 'Ella dijo "Ana" y yo contesté: "Jamás!". Eso fué todo.';
 5  my $a;
 6  ($a = $b) =~ s/".*?!"/-$&-/;
 7  print "$a\n";
 8
 9  $b =~ s/"[^"]*"!"/-$&-/;
10  print "$b\n";
```

Al ejecutar el programa obtenemos:

```
> negynogreedy.pl
Ella dijo -"Ana" y yo contesté: "Jamás!"-. Eso fué todo.
Ella dijo "Ana" y yo contesté: -"Jamás!"-. Eso fué todo.
```

Copia y sustitución simultáneas El operador de *binding* `=~` nos permite “asociar” la variable con la operación de casamiento o sustitución. Si se trata de una sustitución y se quiere conservar la cadena, es necesario hacer una copia:

```
$d = $s;
$d =~ s/esto/por lo otro/;
```

en vez de eso, puedes abreviar un poco usando la siguiente “perla”:

```
($d = $s) =~ s/esto/por lo otro/;
```

Obsérvese la asociación por la izquierda del operador de asignación.

Referencias a Paréntesis Previos

Las referencias relativas permiten escribir expresiones regulares mas reciclables. Véase la documentación en la sección ‘Relative-backreferences’ en `perlretut`:

Counting the opening parentheses to get the correct number for a backreference is error-prone as soon as there is more than one capturing group. A more convenient technique became available with Perl 5.10: relative backreferences. To refer to the immediately preceding capture group one now may write `\g{-1}`, the next but last is available via `\g{-2}`, and so on.

Another good reason in addition to readability and maintainability for using relative backreferences is illustrated by the following example, where a simple pattern for matching peculiar strings is used:

```
1. $a99a = '([a-z])(\d)\2\1'; # matches a11a, g22g, x33x, etc.
```

Now that we have this pattern stored as a handy string, we might feel tempted to use it as a part of some other pattern:

```

1. $line = "code=e99e";
2. if ($line =~ /\^(\\w+)=\$a99a$/){ # unexpected behavior!
3.   print "$1 is valid\n";
4. } else {
5.   print "bad line: '$line'\n";
6. }

```

But this doesn't match – at least not the way one might expect. Only after inserting the interpolated \$a99a and looking at the resulting full text of the regexp is it obvious that the backreferences have backfired – the subexpression (\\w+) has snatched number 1 and demoted the groups in \$a99a by one rank. This can be avoided by using relative backreferences:

```

1. $a99a = '([a-z])(\\d)\\g{-1}\\g{-2}'; # safe for being interpolated

```

El siguiente programa ilustra lo dicho:

```

casiano@millo:~/Lperltesting$ cat -n backreference.pl
1   use strict;
2   use re 'debug';
3
4   my $a99a = '([a-z])(\\d)\\2\\1';
5   my $line = "code=e99e";
6   if ($line =~ /\^(\\w+)=\$a99a$/){ # unexpected behavior!
7     print "$1 is valid\n";
8   } else {
9     print "bad line: '$line'\n";
10  }

```

Sigue la ejecución:

```

casiano@millo:~/Lperltesting$ perl5.10.1 -wd backreference.pl
main::(backreference.pl:4):      my $a99a = '([a-z])(\\d)\\2\\1';
DB<1> c 6
main::(backreference.pl:6):      if ($line =~ /\^(\\w+)=\$a99a$/){ # unexpected behavior!
DB<2> x ($line =~ /\^(\\w+)=\$a99a$/)
empty array
DB<4> $a99a = '([a-z])(\\d)\\g{-1}\\g{-2}'
DB<5> x ($line =~ /\^(\\w+)=\$a99a$/)
0  'code'
1  'e'
2  9

```

Usando Referencias con Nombre (Perl 5.10)

El siguiente texto esta tomado de la sección 'Named-backreferences' en `perlretut`:

Perl 5.10 also introduced named capture buffers and named backreferences. To attach a name to a capturing group, you write either (?<name>...) or (?'name'...). The backreference may then be written as \\g{name}.

It is permissible to attach the same name to more than one group, but then only the leftmost one of the eponymous set can be referenced. Outside of the pattern a named capture buffer is accessible through the %+ hash.

Assuming that we have to match calendar dates which may be given in one of the three formats yyyy-mm-dd, mm/dd/yyyy or dd.mm.yyyy, we can write three suitable patterns where we use 'd', 'm' and 'y' respectively as the names of the buffers capturing the pertaining components of a date. The matching operation combines the three patterns as alternatives:

```

1. $fmt1 = '(?<y>\d\d\d\d)-(?<m>\d\d)-(?<d>\d\d)';
2. $fmt2 = '(?<m>\d\d)/(?<d>\d\d)/(?<y>\d\d\d\d)';
3. $fmt3 = '(?<d>\d\d)\. (?<m>\d\d)\. (?<y>\d\d\d\d)';
4. for my $d qw( 2006-10-21 15.01.2007 10/31/2005 ){
5.     if ( $d =~ m{$fmt1|$fmt2|$fmt3} ){
6.         print "day=${d} month=${m} year=${y}\n";
7.     }
8. }

```

If any of the alternatives matches, the hash %+ is bound to contain the three key-value pairs.

En efecto, al ejecutar el programa:

```

casiano@millo:~/Lperltesting$ cat -n namedbackreferences.pl
1  use v5.10;
2  use strict;
3
4  my $fmt1 = '(?<y>\d\d\d\d)-(?<m>\d\d)-(?<d>\d\d)';
5  my $fmt2 = '(?<m>\d\d)/(?<d>\d\d)/(?<y>\d\d\d\d)';
6  my $fmt3 = '(?<d>\d\d)\. (?<m>\d\d)\. (?<y>\d\d\d\d)';
7
8  for my $d qw( 2006-10-21 15.01.2007 10/31/2005 ){
9      if ( $d =~ m{$fmt1|$fmt2|$fmt3} ){
10         print "day=${d} month=${m} year=${y}\n";
11     }
12 }

```

Obtenemos la salida:

```

casiano@millo:~/Lperltesting$ perl5.10.1 -w namedbackreferences.pl
day=21 month=10 year=2006
day=15 month=01 year=2007
day=31 month=10 year=2005

```

Como se comentó:

... It is permissible to attach the same name to more than one group, but then only the leftmost one of the eponymous set can be referenced.

Veamos un ejemplo:

```

pl@nereida:~/Lperltesting$ perl5.10.1 -wDE 0
main::(-e:1): 0
DB<1> # ... only the leftmost one of the eponymous set can be referenced
DB<2> $r = qr{(?<a>[a-c])(?<a>[a-f])}
DB<3> print ${a} if 'ad' =~ $r
a
DB<4> print ${a} if 'cf' =~ $r
c
DB<5> print ${a} if 'ak' =~ $r

```

Reescribamos el ejemplo de conversión de temperaturas usando paréntesis con nombre:


```

pl@nereida:~/Lperltesting$ cat -n c2f_5_10v2.pl
1  #!/usr/local/bin/perl5_10_1 -w
2  use strict;
3
4  print "Enter a temperature (i.e. 32F, 100C):\n";
5  my $input = <STDIN>;
6  chomp($input);
7
8  $input =~ m/^(
9      (?<fahrenheit>[-+]?[0-9]+(?:\.[0-9]*)?)\s*[fF]
10     |
11     (?<celsius>[-+]?[0-9]+(?:\.[0-9]*)?)\s*[cC]
12     $/x;
13
14  my ($celsius, $fahrenheit);
15  if (exists ${celsius}) {
16      $celsius = ${celsius};
17      $fahrenheit = ($celsius * 9/5)+32;
18  }
19  elsif (exists ${fahrenheit}) {
20      $fahrenheit = ${fahrenheit};
21      $celsius = ($fahrenheit -32)*5/9;
22  }
23  else {
24      die "Expecting a temperature, so don't understand \"$input\".\n";
25  }
26
27  printf "%.2f C = %.2f F\n", $celsius, $fahrenheit;

```

La función `exists` retorna verdadero si existe la clave en el hash y falso en otro caso.

Grupos con Nombre y Factorización

El uso de nombres hace mas robustas y mas factorizables las expresiones regulares. Consideremos la siguiente regexp que usa notación posicional:

```

pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> x "abbacddc" =~ /(.)()\2\1/
0  'a'
1  'b'

```

Supongamos que queremos reutilizar la regexp con repetición

```

DB<2> x "abbacddc" =~ /((.)()\2\1){2}/
empty array

```

¿Que ha ocurrido? La introducción del nuevo paréntesis nos obliga a renombrar las referencias a las posiciones:

```

DB<3> x "abbacddc" =~ /((.)()\3\2){2}/
0  'cddc'
1  'c'
2  'd'
DB<4> "abbacddc" =~ /((.)()\3\2){2}/; print "$&\n"
abbacddc

```

Esto no ocurre si utilizamos nombres. El operador `\k<a>` sirve para hacer referencia al valor que ha casado con el paréntesis con nombre `a`:

```
DB<5> x "abbacddc" =~ /((?<a>.) (?<b>.) \k<b>\k<a>){2}/
0 'cddc'
1 'c'
2 'd'
```

El uso de grupos con nombre y `\k1` en lugar de referencias numéricas absolutas hace que la regexp sea mas reutilizable.

LLlamadas a expresiones regulares via paréntesis con memoria

Es posible también llamar a la expresión regular asociada con un paréntesis.

Este parrafo tomado de la sección 'Extended-Patterns' en `perlre` explica el modo de uso:

`(?PARNO) (?-PARNO) (?+PARNO) (?R) (?0)`

PARNO is a sequence of digits (not starting with 0) whose value reflects the paren-number of the capture buffer to recurse to.

....

Capture buffers contained by the pattern will have the value as determined by the outermost recursion.

If PARNO is preceded by a plus or minus sign then it is assumed to be relative, with negative numbers indicating preceding capture buffers and positive ones following. Thus (?-1) refers to the most recently declared buffer, and (?+1) indicates the next buffer to be declared.

Note that the counting for relative recursion differs from that of relative backreferences, in that with recursion unclosed buffers are included.

Veamos un ejemplo:

```
casiano@millo:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> x "AABB" =~ /(A)(?-1)(?+1)(B)/
0 'A'
1 'B'
# Parenthesis:      1    2    2                      1
DB<2> x 'ababa' =~ / ^((?:([ab])(?1)\g{-1}|[ab]?))$/
0 'ababa'
1 'a'
DB<3> x 'bbabababb' =~ / ^((?:([ab])(?1)\g{-1}|[ab]?))$/
0 'bbabababb'
1 'b'
```

Véase también:

- Perl Training Australia: Regular expressions in Perl 5.10
- Perl 5.10 Advanced Regular Expressions by Yves Orton
- Gabor: Regular Expressions in Perl 5.10

¹ Una diferencia entre `\k` y `\g` es que el primero sólo admite un nombre como argumento mientras que `\g` admite enteros

Reutilizando Expresiones Regulares

La siguiente reescritura de nuestro ejemplo básico utiliza el módulo `Regexp::Common` para factorizar la expresión regular:

```
casiano@millor:~/src/perl/perltesting$ cat -n c2f_5_10v3.pl
1  #!/soft/perl5lib/bin/perl5.10.1 -w
2  use strict;
3  use Regexp::Common;
4
5  print "Enter a temperature (i.e. 32F, 100C):\n";
6  my $input = <STDIN>;
7  chomp($input);
8
9  $input =~ m/~
10      (?<fahrenheit>$RE{num}{real})\s*[fF]
11      |
12      (?<celsius>$RE{num}{real})\s*[cC]
13      $/x;
14
15  my ($celsius, $fahrenheit);
16  if ('celsius' ~~ %+) {
17      $celsius = ${+{celsius}};
18      $fahrenheit = ($celsius * 9/5)+32;
19  }
20  elsif ('fahrenheit' ~~ %+) {
21      $fahrenheit = ${+{fahrenheit}};
22      $celsius = ($fahrenheit -32)*5/9;
23  }
24  else {
25      die "Expecting a temperature, so don't understand \"$input\".\n";
26  }
27
28  printf "%.2f C = %.2f F\n", $celsius, $fahrenheit;
```

Véase:

- La documentación del módulo `Regexp::Common` por Abigail
- Smart Matching: Perl Training Australia: Smart Match
- Rafael García Suárez: la sección 'Smart-matching-in-detail' en `perlsyn`
- Enrique Nell (Barcelona Perl Mongers): Novedades en Perl 5.10

El Módulo `Regexp::Common`

El módulo `Regexp::Common` provee un extenso número de expresiones regulares que son accesibles vía el hash `%RE`. sigue un ejemplo de uso:

```
casiano@millor:~/Lperltesting$ cat -n regexpcommonsynopsis.pl
1  use strict;
2  use Perl6::Say;
3  use Regexp::Common;
4
5  while (<>) {
6      say q{a number}          if /$RE{num}{real}/;
7  }
```

```

8      say q{a ['"'] quoted string} if /$RE{quoted}/;
9
10     say q{a /.../ sequence}      if m{$RE{delimited}{'-delim'=>'/'}};
11
12     say q{balanced parentheses}  if /$RE{balanced}{'-parens'=>'()'}/;
13
14     die q{a #*%-ing word}."\\n"  if /$RE{profanity}/;
15
16 }
17

```

Sigue un ejemplo de ejecución:

```

casiano@millo:~/Lperltesting$ perl regexprcommonsynopsis.pl
43
a number
"2+2 es" 4
a number
a ['"'] quoted string
x/y/z
a /.../ sequence
(2*(4+5/(3-2)))
a number
balanced parentheses
fuck you!
a #*%-ing word

```

El siguiente fragmento de la documentación de Regexp::Common explica el modo simplificado de uso:

To access a particular pattern, %RE is treated as a hierarchical hash of hashes (of hashes...), with each successive key being an identifier. For example, to access the pattern that matches real numbers, you specify:

```
$RE{num}{real}
```

and to access the pattern that matches integers:

```
$RE{num}{int}
```

Deeper layers of the hash are used to specify flags: arguments that modify the resulting pattern in some way.

- *The keys used to access these layers are prefixed with a minus sign and may have a value;*
- *if a value is given, it's done by using a multidimensional key.*

For example, to access the pattern that matches base-2 real numbers with embedded commas separating groups of three digits (e.g. 10,101,110.110101101):

```
$RE{num}{real}{-base => 2}{-sep => ', '}{-group => 3}
```

Through the magic of Perl, these flag layers may be specified in any order (and even interspersed through the identifier keys!) so you could get the same pattern with:

`$RE{num}{real}{-sep => ', '}{-group => 3}{-base => 2}`

or:

`$RE{num}{-base => 2}{real}{-group => 3}{-sep => ', '}`

or even:

`$RE{-base => 2}{-group => 3}{-sep => ', '}{num}{real}`

etc.

Note, however, that the relative order of amongst the identifier keys is significant. That is:

`$RE{list}{set}`

would not be the same as:

`$RE{set}{list}`

Veamos un ejemplo con el depurador:

```
casiano@millo:~/Lperltesting$ perl -MRegex::Common -wde 0
main::(-e:1): 0
DB<1> x 'numero: 10,101,110.110101101 101.1e-1 234' =~ m{($RE{num}{real}{-base => 2}{-sep =>
0 '10,101,110.110101101'
1 '101.1e-1'
```

La expresión regular para un número real es relativamente compleja:

```
casiano@millo:~/src/perl/perltesting$ perl5.10.1 -wd c2f_5_10v3.pl
main::(c2f_5_10v3.pl:5): print "Enter a temperature (i.e. 32F, 100C):\n";
DB<1> p $RE{num}{real}
(?: (?:i) (?: [+ -] ?) (?: (?: [0123456789] | [.] ) (?: [0123456789] *) (?: (?: [.] ) (?: [0123456789] {0,} )) ?) (?: (?
```

Si se usa la opción `-keep` el patrón proveído usa paréntesis con memoria:

```
casiano@millo:~/Lperltesting$ perl -MRegex::Common -wde 0
main::(-e:1): 0
DB<2> x 'one, two, three, four, five' =~ /$RE{list}{-pat => '\w+'}/
0 1
DB<3> x 'one, two, three, four, five' =~ /$RE{list}{-pat => '\w+'}{-keep}/
0 'one, two, three, four, five'
1 ', '
```

Smart Matching

Perl 5.10 introduce el operador de smart matching. El siguiente texto es tomado casi verbatim del site de la compañía Perl Training Australia²:

² This Perl tip and associated text is copyright Perl Training Australia

Perl 5.10 introduces a new-operator, called *smart-match*, written `~~`. As the name suggests, *smart-match* tries to compare its arguments in an intelligent fashion. Using *smart-match* effectively allows many complex operations to be reduced to very simple statements.

Unlike many of the other features introduced in Perl 5.10, there's no need to use the *feature* pragma to enable *smart-match*, as long as you're using 5.10 it's available.

The *smart-match* operator is always commutative. That means that `$x ~~ $y` works the same way as `$y ~~ $x`. You'll never have to remember which order to place to your operands with *smart-match*. *Smart-match in action.*

As a simple introduction, we can use *smart-match* to do a simple string comparison between simple scalars. For example:

```
use feature qw(say);

my $x = "foo";
my $y = "bar";
my $z = "foo";

say '$x and $y are identical strings' if $x ~~ $y;
say '$x and $z are identical strings' if $x ~~ $z;    # Printed
```

If one of our arguments is a number, then a numeric comparison is performed:

```
my $num    = 100;
my $input  = <STDIN>;

say 'You entered 100' if $num ~~ $input;
```

This will print our message if our user enters 100, 100.00, +100, 1e2, or any other string that looks like the number 100.

We can also smart-match against a regexp:

```
my $input  = <STDIN>;

say 'You said the secret word!' if $input ~~ /xyzzzy/;
```

Smart-matching with a regexp also works with saved regexps created with `qr`.

*So we can use *smart-match* to act like `eq`, `==` and `=~`, so what? Well, it does much more than that.*

*We can use *smart-match* to search a list:*

```
casiano@millo:~/Lperltesting$ perl5.10.1 -wdE 0
main::(-e:1):    0
DB<1> @friends = qw(Frodo Meriadoc Pippin Samwise Gandalf)
DB<2> print "You're a friend" if 'Pippin' ~~ @friends
You're a friend
DB<3> print "You're a friend" if 'Mordok' ~~ @friends
```

*It's important to note that searching an array with *smart-match* is extremely fast. It's faster than using `grep`, it's faster than using `first` from `Scalar::Util`, and it's faster than walking through the loop with `foreach`, even if you do know all the clever optimisations.*

Esta es la forma típica de buscar un elemento en un array en versiones anteriores a la 5.10:

```
casiano@millo:~$ perl -wde 0
main::(-e:1): 0
DB<1> use List::Util qw{first}
DB<2> @friends = qw(Frodo Meriadoc Pippin Samwise Gandalf)
DB<3> x first { $_ eq 'Pippin' } @friends
0 'Pippin'
DB<4> x first { $_ eq 'Mordok' } @friends
0 undef
```

We can also use smart-match to compare arrays:

```
DB<4> @foo = qw(x y z xyzzy ninja)
DB<5> @bar = qw(x y z xyzzy ninja)
DB<7> print "Identical arrays" if @foo ~~ @bar
Identical arrays
DB<8> @bar = qw(x y z xyzzy nOnjA)
DB<9> print "Identical arrays" if @foo ~~ @bar
DB<10>
```

And even search inside an array using a string:

```
DB<11> x @foo = qw(x y z xyzzy ninja)
0 'x'
1 'y'
2 'z'
3 'xyzzy'
4 'ninja'
DB<12> print "Array contains a ninja " if @foo ~~ 'ninja'
```

or using a regexp:

```
DB<13> print "Array contains magic pattern" if @foo ~~ /xyz/
Array contains magic pattern
DB<14> print "Array contains magic pattern" if @foo ~~ /\d+/'
```

Smart-match works with array references, too³:

```
DB<16> $array_ref = [ 1..10 ]
DB<17> print "Array contains 10" if 10 ~~ $array_ref
Array contains 10
DB<18> print "Array contains 10" if $array_ref ~~ 10
DB<19>
```

En el caso de un número y un array devuelve cierto si el escalar aparece en un array anidado:

```
casiano@millo:~/Lperltesting$ perl5.10.1 -E 'say "ok" if 42 ~~ [23, 17, [40..50], 70];'
ok
casiano@millo:~/Lperltesting$ perl5.10.1 -E 'say "ok" if 42 ~~ [23, 17, [50..60], 70];'
casiano@millo:~/Lperltesting$
```

Of course, we can use smart-match with more than just arrays and scalars, it works with searching for the key in a hash, too!

³En este caso la conmutatividad no funciona

```

DB<19> %colour = ( sky    => 'blue', grass => 'green', apple => 'red',)
DB<20> print "I know the colour" if 'grass' ~~ %colour
I know the colour
DB<21> print "I know the colour" if 'cloud' ~~ %colour
DB<22>
DB<23> print "A key starts with 'gr'" if %colour ~~ /^gr/
A key starts with 'gr'
DB<24> print "A key starts with 'clou'" if %colour ~~ /^clou/
DB<25>

```

You can even use it to see if the two hashes have identical keys:

```

DB<26> print 'Hashes have identical keys' if %taste ~~ %colour;
Hashes have identical keys

```

La conducta del operador de smart matching viene dada por la siguiente tabla tomada de la sección 'Smart-matching-in-detail' en `perlsyn`:

The behaviour of a smart match depends on what type of thing its arguments are. The behaviour is determined by the following table: the first row that applies determines the match behaviour (which is thus mostly determined by the type of the right operand). Note that the smart match implicitly dereferences any non-blessed hash or array ref, so the `"Hash.and.Array.entries` apply in those cases. (For blessed references, the `.Object.entries` apply.)

Note that the "Matching Code" column is not always an exact rendition. For example, the smart match operator short-circuits whenever possible, but `grep` does not.

\$a	\$b	Type of Match Implied	Matching Code
=====	=====	=====	=====
Any	undef	undefined	!defined \$a
Any	Object	invokes <code>~~</code> overloading on \$object, or dies	
Hash	CodeRef	sub truth for each key[1]	!grep { !\$b->(\$_) } keys %\$a
Array	CodeRef	sub truth for each elt[1]	!grep { !\$b->(\$_) } @\$a
Any	CodeRef	scalar sub truth	\$b->(\$a)
Hash	Hash	hash keys identical (every key is found in both hashes)	
Array	Hash	hash slice existence	grep { exists \$b->{\$_} } @\$a
Regex	Hash	hash key grep	grep /\$a/, keys %\$b
undef	Hash	always false (undef can't be a key)	
Any	Hash	hash entry existence	exists \$b->{\$a}
Hash	Array	hash slice existence	grep { exists \$a->{\$_} } @\$b
Array	Array	arrays are comparable[2]	
Regex	Array	array grep	grep /\$a/, @\$b
undef	Array	array contains undef	grep !defined, @\$b
Any	Array	match against an array element[3]	grep \$a ~~ \$_, @\$b
Hash	Regex	hash key grep	grep /\$b/, keys %\$a
Array	Regex	array grep	grep /\$b/, @\$a
Any	Regex	pattern match	\$a =~ /\$b/

Object	Any	invokes ~~ overloading on \$object, or falls back:	
Any	Num	numeric equality	\$a == \$b
Num	numish[4]	numeric equality	\$a == \$b
undef	Any	undefined	!defined(\$b)
Any	Any	string equality	\$a eq \$b

Ejercicios

Ejercicio 3.1.4. ■ *Indique la salida del siguiente programa:*

```

1  pl@nereida:~/Lperltesting$ cat twonumbers.pl
2  $_ = "I have 2 numbers: 53147";
3  @pats = qw{
4      (.*)(\d*)
5      (.*)(\d+)
6      (.*?)(\d*)
7      (.*?)(\d+)
8      (.*)(\d+)$
9      (.*?)(\d+)$
10     (.*)\b(\d+)$
11     (.*\D)(\d+)$
12 };
13
14 print "$_\n";
15 for $pat (@pats) {
16     printf "%-12s ", $pat;
17     <>;
18     if ( /$pat/ ) {
19         print "<$1> <$2>\n";
20     } else {
21         print "FAIL\n";
22     }
23 }
```

3.1.2. Depuración de Expresiones Regulares

Para obtener información sobre la forma en que es compilada una expresión regular y como se produce el proceso de matching podemos usar la opción 'debug' del módulo `re`. La versión de Perl 5.10 da una información algo mas legible que la de las versiones anteriores:

```
pl@nereida:~/Lperltesting$ perl5_10_1 -wde 0
```

```
Loading DB routines from perl5db.pl version 1.32
Editor support available.
```

```
Enter h or 'h h' for help, or 'man perldebug' for more help.
```

```
main::(-e:1): 0
DB<1> use re 'debug'; 'astr' =~ m{[sf].r}
Compiling REx "[sf].r"
Final program:
  1: ANYOF[fs] [] (12)
 12: REG_ANY (13)
 13: EXACT <r> (15)
```

```

15: END (0)
anchored "r" at 2 (checking anchored) stclass ANYOF[fs][] minlen 3
Guessing start of match in sv for REx "[sf].r" against "astr"
Found anchored substr "r" at offset 3...
Starting position does not contradict /^/m...
start_shift: 2 check_at: 3 s: 1 endpos: 2
Does not contradict STCLASS...
Guessed: match at offset 1
Matching REx "[sf].r" against "str"
  1 <a> <str>          | 1:ANYOF[fs][] (12)
  2 <as> <tr>          | 12:REG_ANY(13)
  3 <ast> <r>          | 13:EXACT <r>(15)
  4 <astr> <>         | 15:END(0)
Match successful!
Freeing REx: "[sf].r"

```

Si se usa la opción `debug` de `re` con objetos expresión regular, se obtendrá información durante el proceso de matching:

```

DB<3> use re 'debug'; $re = qr{[sf].r}
Compiling REx "[sf].r"
Final program:
  1: ANYOF[fs][] (12)
 12: REG_ANY (13)
 13: EXACT <r> (15)
 15: END (0)
anchored "r" at 2 (checking anchored) stclass ANYOF[fs][] minlen 3

DB<4> 'astr' =~ $re
Guessing start of match in sv for REx "[sf].r" against "astr"
Found anchored substr "r" at offset 3...
Starting position does not contradict /^/m...
start_shift: 2 check_at: 3 s: 1 endpos: 2
Does not contradict STCLASS...
Guessed: match at offset 1
Matching REx "[sf].r" against "str"
  1 <a> <str>          | 1:ANYOF[fs][] (12)
  2 <as> <tr>          | 12:REG_ANY(13)
  3 <ast> <r>          | 13:EXACT <r>(15)
  4 <astr> <>         | 15:END(0)
Match successful!

```

3.1.3. Tablas de Escapes, Metacaracteres, Cuantificadores, Clases

Sigue una sección de tablas con notaciones tomada de `perlre`:

Metacharacters

The following metacharacters have their standard `egrep`-ish meanings:

1. \ Quote the next metacharacter
2. ^ Match the beginning of the line
3. . Match any character (except newline)
4. \$ Match the end of the line (or before newline at the end)
5. | Alternation

6. () Grouping
7. [] Character class

Standard greedy quantifiers

The following standard greedy quantifiers are recognized:

1. * Match 0 or more times
2. + Match 1 or more times
3. ? Match 1 or 0 times
4. {n} Match exactly n times
5. {n,} Match at least n times
6. {n,m} Match at least n but not more than m times

Non greedy quantifiers

The following non greedy quantifiers are recognized:

1. *? Match 0 or more times, not greedily
2. +? Match 1 or more times, not greedily
3. ?? Match 0 or 1 time, not greedily
4. {n}? Match exactly n times, not greedily
5. {n,}? Match at least n times, not greedily
6. {n,m}? Match at least n but not more than m times, not greedily

Possesive quantifiers

The following possesive quantifiers are recognized:

1. ++ Match 0 or more times and give nothing back
2. ++ Match 1 or more times and give nothing back
3. ?+ Match 0 or 1 time and give nothing back
4. {n}+ Match exactly n times and give nothing back (redundant)
5. {n,}+ Match at least n times and give nothing back
6. {n,m}+ Match at least n but not more than m times and give nothing back

Escape sequences

1. \t tab (HT, TAB)
2. \n newline (LF, NL)
3. \r return (CR)
4. \f form feed (FF)
5. \a alarm (bell) (BEL)
6. \e escape (think troff) (ESC)
7. \033 octal char (example: ESC)
8. \x1B hex char (example: ESC)
9. \x{263a} long hex char (example: Unicode SMILEY)
10. \cK control char (example: VT)
11. \N{name} named Unicode character
12. \l lowercase next char (think vi)
13. \u uppercase next char (think vi)
14. \L lowercase till \E (think vi)
15. \U uppercase till \E (think vi)
16. \E end case modification (think vi)
17. \Q quote (disable) pattern metacharacters till \E

Ejercicio 3.1.5. *Explique la salida:*

```
casiano@tonga:~$ perl -wde 0
main::(-e:1): 0
DB<1> $x = '([a-z]+)'
DB<2> x 'hola' =~ /$x/
0 'hola'
DB<3> x 'hola' =~ /\Q$x/
empty array
DB<4> x '([a-z]+)' =~ /\Q$x/
0 1
```

Character Classes and other Special Escapes

1. \w Match a "word" character (alphanumeric plus "_")
2. \W Match a non-"word" character
3. \s Match a whitespace character
4. \S Match a non-whitespace character
5. \d Match a digit character
6. \D Match a non-digit character
7. \pP Match P, named property. Use \p{Prop} for longer names.
8. \PP Match non-P
9. \X Match eXtended Unicode "combining character sequence",
equivalent to (?>\PM\pM*)
11. \C Match a single C char (octet) even under Unicode.
12. NOTE: breaks up characters into their UTF-8 bytes,
13. so you may end up with malformed pieces of UTF-8.
14. Unsupported in lookbehind.
15. \1 Backreference to a specific group.
16. '1' may actually be any positive integer.
17. \g1 Backreference to a specific or previous group,
18. \g{-1} number may be negative indicating a previous buffer and may
19. optionally be wrapped in curly brackets for safer parsing.
20. \g{name} Named backreference
21. \k<name> Named backreference
22. \K Keep the stuff left of the \K, don't include it in \$&
23. \v Vertical whitespace
24. \V Not vertical whitespace
25. \h Horizontal whitespace
26. \H Not horizontal whitespace
27. \R Linebreak

Zero width assertions

Perl defines the following zero-width assertions:

1. \b Match a word boundary
2. \B Match except at a word boundary
3. \A Match only at beginning of string
4. \Z Match only at end of string, or before newline at the end
5. \z Match only at end of string
6. \G Match only at pos() (e.g. at the end-of-match position
7. of prior m//g)

The POSIX character class syntax

The POSIX character class syntax:

1. `[:class:]`

is also available. Note that the `[` and `]` brackets are literal; they must always be used within a character class expression.

1. `# this is correct:`
2. `$string =~ /[[:alpha:]]/;`
- 3.
4. `# this is not, and will generate a warning:`
5. `$string =~ /[:alpha:]/;`

Available classes

The available classes and their backslash equivalents (if available) are as follows:

1. `alpha`
2. `alnum`
3. `ascii`
4. `blank`
5. `cntrl`
6. `digit \d`
7. `graph`
8. `lower`
9. `print`
10. `punct`
11. `space \s`
12. `upper`
13. `word \w`
14. `xdigit`

For example use `[:upper:]` to match all the uppercase characters. Note that the `[]` are part of the `[::]` construct, not part of the whole character class. For example:

1. `[01[:alpha:]]%`

matches zero, one, any alphabetic character, and the percent sign.

Equivalences to Unicode

The following equivalences to Unicode `\p{}` constructs and equivalent backslash character classes (if available), will hold:

1. `[[:...:]] \p{...} backslash`
- 2.
3. `alpha IsAlpha`
4. `alnum IsAlnum`
5. `ascii IsASCII`
6. `blank`
7. `cntrl IsCntrl`
8. `digit IsDigit \d`
9. `graph IsGraph`
10. `lower IsLower`
11. `print IsPrint`
12. `punct IsPunct`
13. `space IsSpace`
14. `IsSpacePerl \s`
15. `upper IsUpper`
16. `word IsWord \w`
17. `xdigit IsXDigit`

Negated character classes

You can negate the `[::]` character classes by prefixing the class name with a `'^'`. This is a Perl extension. For example:

1. POSIX traditional Unicode
- 2.
3. `[[:^digit:]] \D \P{IsDigit}`
4. `[[:^space:]] \S \P{IsSpace}`
5. `[[:^word:]] \W \P{IsWord}`

3.1.4. Variables especiales después de un emparejamiento

Después de un emparejamiento con éxito, las siguientes variables especiales quedan definidas:

<code>\$&</code>	El texto que casó
<code>\$'</code>	El texto que está a la izquierda de lo que casó
<code>\$'</code>	El texto que está a la derecha de lo que casó
<code>\$1, \$2, \$3, etc.</code>	Los textos capturados por los paréntesis
<code>\$+</code>	Una copia del <code>\$1, \$2, ...</code> con número mas alto
<code>@-</code>	Desplazamientos de las subcadenas que casan en <code>\$1 ...</code>
<code>@+</code>	Desplazamientos de los finales de las subcadenas en <code>\$1 ...</code>
<code>\$#-</code>	El índice del último paréntesis que casó
<code>\$#+</code>	El índice del último paréntesis en la última expresión regular

Las Variables de match, pre-match y post-match

Ejemplo:

```
1 #!/usr/bin/perl -w
2 if ("Hello there, neighbor" =~ /\s(\w+),/) {
3   print "That was: ($')($&($').\n",
4 }
```

```
> matchvariables.pl
```

```
That was: (Hello)( there,)( neighbor).
```

El uso de estas variables tenía un efecto negativo en el rendimiento de la regexp. Véase `perlfaq6` la sección `Why does using $&, $', or $' slow my program down?`.

*Once Perl sees that you need one of these variables anywhere in the program, it provides them on each and every pattern match. That means that on every pattern match the entire string will be copied, part of it to `$'`, part to `$&`, and part to `$'`. Thus the penalty is most severe with long strings and patterns that match often. Avoid `$&`, `$'`, and `$'` if you can, but if you can't, once you've used them at all, use them at will because you've already paid the price. Remember that some algorithms really appreciate them. As of the 5.005 release, the `$&` variable is no longer *expensive* the way the other two are.*

Since Perl 5.6.1 the special variables `@-` and `@+` can functionally replace `$'`, `$&` and `$'`. These arrays contain pointers to the beginning and end of each match (see `perlvar` for the full story), so they give you essentially the same information, but without the risk of excessive string copying.

Perl 5.10 added three specials, `${^MATCH}`, `${^PREMATCH}`, and `${^POSTMATCH}` to do the same job but without the global performance penalty. Perl 5.10 only sets these variables if you compile or execute the regular expression with the `/p` modifier.

```

pl@nereida:~/Lperltesting$ cat ampersandoldway.pl
#!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w
use strict;
use Benchmark qw(cmpthese timethese);

'hola juan' =~ /ju/;
my ($a, $b, $c) = ($', $&, $');

cmpthese( -1, {
    oldway => sub { 'hola juan' =~ /ju/ },
});
pl@nereida:~/Lperltesting$ cat ampersandnewway.pl
#!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w
use strict;
use Benchmark qw(cmpthese timethese);

'hola juan' =~ /ju/p;
my ($a, $b, $c) = (${'^PREMATCH'}, ${'^MATCH'}, ${'^POSTMATCH'});

cmpthese( -1, {
    newway => sub { 'hola juan' =~ /ju/ },
});

pl@nereida:~/Lperltesting$ time ./ampersandoldway.pl
Rate oldway
oldway 2991861/s      --

real    0m3.761s
user    0m3.740s
sys      0m0.020s
pl@nereida:~/Lperltesting$ time ./ampersandnewway.pl
Rate newway
newway 8191999/s      --

real    0m6.721s
user    0m6.704s
sys      0m0.016s

```

Véase

- perlvar (busque por \$MATCH)

Texto Asociado con el Último Paréntesis

La variable \$+ contiene el texto que casó con el último paréntesis en el patrón. Esto es útil en situaciones en las cuáles una de un conjunto de alternativas casa, pero no sabemos cuál:

```

DB<9> "Revision: 4.5" =~ /Version: (.*)|Revision: (.*)/ && ($rev = $+);
DB<10> x $rev
0 4.5
DB<11> "Version: 4.5" =~ /Version: (.*)|Revision: (.*)/ && ($rev = $+);
DB<12> x $rev
0 4.5

```

Los Offsets de los Inicios de los Casamientos: @-

El vector @- contiene los *offsets* o desplazamientos de los casamientos en la última expresión regular. La entrada \$-[0] es el desplazamiento del último casamiento con éxito y \$-[n] es el desplazamiento de la subcadena que casa con el n-ésimo paréntesis (o undef si el paréntesis no casó). Por ejemplo:

```
# 012345678
DB<1> $z = "hola13.47"
DB<2> if ($z =~ m{a(\d+)(\.(\\d+))?}) { print "@-\n"; }
3 4 6 7
```

El resultado se interpreta como sigue:

- 3 = desplazamiento de comienzo de \$& = a13.47
- 4 = desplazamiento de comienzo de \$1 = 13
- 6 = desplazamiento de comienzo de \$2 = .
- 7 = desplazamiento de comienzo de \$3 = 47

Esto es lo que dice perlvar sobre @-:

This array holds the offsets of the beginnings of the last successful submatches in the currently active dynamic scope. \$-[0] is the offset into the string of the beginning of the entire match. The nth element of this array holds the offset of the nth submatch, so \$-[1] is the offset where \$1 begins, \$-[2] the offset where \$2 begins, and so on.

After a match against some variable \$var:

```
$' is the same as substr($var, 0, $-[0])
$& is the same as substr($var, $-[0], $+[0] - $-[0])
$' is the same as substr($var, $+[0])
$1 is the same as substr($var, $-[1], $+[1] - $-[1])
$2 is the same as substr($var, $-[2], $+[2] - $-[2])
$3 is the same as substr($var, $-[3], $+[3] - $-[3])
```

Desplazamientos de los Finales de los Emparejamientos: @+

El array @+ contiene los desplazamientos de los finales de los emparejamientos. La entrada \$+[0] contiene el desplazamiento del final de la cadena del emparejamiento completo. Siguiendo con el ejemplo anterior:

```
# 0123456789
DB<17> $z = "hola13.47x"
DB<18> if ($z =~ m{a(\d+)(\.(\\d+))?}) { print "@+\n"; }
9 6 7 9
```

El resultado se interpreta como sigue:

- 9 = desplazamiento final de \$& = a13.47x
- 6 = desplazamiento final de \$1 = 13
- 7 = desplazamiento final de \$2 = .
- 9 = desplazamiento final de \$3 = 47

Número de paréntesis en la última regexp con éxito

Se puede usar `$$+` para determinar cuantos parentesis había en el último emparejamiento que tuvo éxito.

```
DB<29> $z = "h"
DB<30> print "$#+\n" if ($z =~ m{(a)(b)}) || ($z =~ m{(h)(.)?(.)?})
3
DB<31> $z = "ab"
DB<32> print "$#+\n" if ($z =~ m{(a)(b)}) || ($z =~ m{(h)(.)?(.)?})
2
```

Índice del Último Paréntesis

La variable `$$-` contiene el índice del último paréntesis que casó. Observe la siguiente ejecución con el depurador:

```
DB<1> $x = '13.47'; $y = '125'
DB<2> if ($y =~ m{(\d+)(\.(d+))?}) { print "last par = $$-, content = $+\n"; }
last par = 1, content = 125
DB<3> if ($x =~ m{(\d+)(\.(d+))?}) { print "last par = $$-, content = $+\n"; }
last par = 3, content = 47
```

@- y @+ no tienen que tener el mismo tamaño

En general no puede asumirse que `@-` y `@+` sean del mismo tamaño.

```
DB<1> "a" =~ /(a)|(b)/; @a = @-; @b = @+
DB<2> x @a
0 0
1 0
DB<3> x @b
0 1
1 1
2 undef
```

Véase También

Para saber más sobre las variables especiales disponibles consulte

- `perldoc perlretut`
- `perldoc perlvar`.

3.1.5. Ambito Automático

Como sabemos, ciertas variables (como `$1`, `$&` ...) reciben automáticamente un valor con cada operación de “matching”.

Considere el siguiente código:

```
if (m/(...)/) {
    &do_something();
    print "the matched variable was $1.\n";
}
```

Puesto que `$1` es automáticamente declarada `local` a la entrada de cada bloque, no importa lo que se haya hecho en la función `&do_something()`, el valor de `$1` en la sentencia `print` es el correspondiente al “matching” realizado en el `if`.

3.1.6. Opciones

Modificador	Significado
e	evaluar: evaluar el lado derecho de una sustitución como una expresión
g	global: Encontrar todas las ocurrencias
i	ignorar: no distinguir entre mayúsculas y minúsculas
m	multilínea (^ y \$ casan con \n internos)
o	optimizar: compilar una sola vez
s	^ y \$ ignoran \n pero el punto . “casa” con \n
x	extendida: permitir comentarios

El Modificador /g La conducta de este modificador depende del contexto. En un contexto de listas devuelve una lista con todas las subcadenas casadas por todos los paréntesis en la expresión regular. Si no hubieran paréntesis devuelve una lista con todas las cadenas casadas (como si hubiera paréntesis alrededor del patrón global).

```
1 #!/usr/bin/perl -w
2 ($one, $five, $fifteen) = ('uptime' =~ /(\d+\.\d+)/g);
3 print "$one, $five, $fifteen\n";
```

Observe la salida:

```
> uptime
1:35pm up 19:22, 0 users, load average: 0.01, 0.03, 0.00
> glist.pl
0.01, 0.03, 0.00
```

En un contexto escalar `m//g` itera sobre la cadena, devolviendo cierto cada vez que casa, y falso cuando deja de casar. En otras palabras, recuerda donde se quedo la última vez y se recomienza la búsqueda desde ese punto. Se puede averiguar la posición del emparejamiento utilizando la función `pos`. Si por alguna razón modificas la cadena en cuestión, la posición de emparejamiento se reestablece al comienzo de la cadena.

```
1 #!/usr/bin/perl -w
2 # count sentences in a document
3 #defined as ending in [.!?] perhaps with
4 # quotes or parens on either side.
5 $/ = ""; # paragraph mode
6 while ($paragraph = <>) {
7     print $paragraph;
8     while ($paragraph =~ /[a-z](['"])*[.!?]+(['"])*\s/g) {
9         $sentences++;
10    }
11 }
12 print "$sentences\n";
```

Observe el uso de la variable especial `$/`. Esta variable contiene el separador de registros en el fichero de entrada. Si se iguala a la cadena vacía usará las líneas en blanco como separadores. Se le puede dar el valor de una cadena multicarácter para usarla como delimitador. Nótese que establecerla a `\n\n` es diferente de asignarla a `"`. Si se deja `undef`, la siguiente lectura leerá todo el fichero.

Sigue un ejemplo de ejecución. El programa se llama `gscalar.pl`. Introducimos el texto desde STDIN. El programa escribe el número de párrafos:

```
> gscalar.pl
este primer parrafo. Sera seguido de un
segundo parrafo.
```

"Cita de Seneca".

3

La opción e: Evaluación del remplazo La opción /e permite la evaluación como expresión perl de la cadena de reemplazo (En vez de considerarla como una cadena delimitada por doble comilla).

```
1 #!/usr/bin/perl -w
2 $_ = "abc123xyz\n";
3 s/\d+/$&*2/e;
4 print;
5 s/\d+/sprintf("%5d",$&)/e;
6 print;
7 s/\w/$& x 2/eg;
8 print;
```

El resultado de la ejecución es:

```
> replacement.pl
abc246xyz
abc 246xyz
aabbcc 224466xxyyzz
```

Véase un ejemplo con anidamiento de /e:

```
1 #!/usr/bin/perl
2 $a = "one";
3 $b = "two";
4 $_ = '$a $b';
5 print "_ = $_\n\n";
6 s/(\$(w+))/$_/ge;
7 print "After 's/(\$(w+))/$_/ge' _ = $_\n\n";
8 s/(\$(w+))/$_/gee;
9 print "After 's/(\$(w+))/$_/gee' _ = $_\n\n";
```

El resultado de la ejecución es:

```
> enested.pl
_ = $a $b
```

After 's/(\$w+)/\$b/ge' _ = \$a \$b

After 's/(\$w+)/\$b/gee' _ = one two

He aquí una solución que hace uso de e al siguiente ejercicio (véase 'Regex to add space after punctuation sign' en PerlMonks) Se quiere poner un espacio en blanco después de la aparición de cada coma:

```
s/,/, /g;
```

pero se quiere que la sustitución no tenga lugar si la coma esta incrustada entre dos dígitos. Además se pide que si hay ya un espacio después de la coma, no se duplique

```
s/(\d[,.\d]|(?!\s))/ $1 || "$2 "/ge;
```

Se hace uso de un lookahead negativo (?!\s). Véase la sección 3.2.3 para entender como funciona un lookahead negativo.

3.2. Algunas Extensiones

3.2.1. Comentarios

(*?#text*) Un comentario. Se ignora *text*. Si se usa la opción *x* basta con poner *#*.

3.2.2. Modificadores locales

Los modificadores de la conducta de una expresión regular pueden ser empotrados en una subexpresión usando el formato (*?pimsx-imsx*).

Véase el correspondiente texto *Extended Patterns* de la sección 'Extended-Patterns' en *perlre*:

One or more embedded pattern-match modifiers, to be turned on (or turned off, if preceded by '-') for the remainder of the pattern or the remainder of the enclosing pattern group (if any). This is particularly useful for dynamic patterns, such as those read in from a configuration file, taken from an argument, or specified in a table somewhere. Consider the case where some patterns want to be case sensitive and some do not: The case insensitive ones merely need to include (?i) at the front of the pattern. For example:

```
1. $pattern = "foobar";
2. if ( /$pattern/i ) { }
3.
4. # more flexible:
5.
6. $pattern = "(?i)foobar";
7. if ( /$pattern/ ) { }
```

These modifiers are restored at the end of the enclosing group. For example,

```
1. ( (?i) blah ) \s+ \1
```

will match blah in any case, some spaces, and an exact (including the case!) repetition of the previous word, assuming the /x modifier, and no /i modifier outside this group.

El siguiente ejemplo extiende el ejemplo visto en la sección 3.1.1 eliminando los comentarios */* ... */* y *// ...* de un programa C. En dicho ejemplo se usaba el modificador *s* para hacer que el punto casara con cualquier carácter:

```
casiano@tonga:~/Lperltesting$ cat -n extendedcomments.pl
1  #!/usr/bin/perl -w
2  use strict;
3
4  my $programe = shift @ARGV or die "Usage:\n$0 prog.c\n";
5  open(my $PROGRAM,"<$programe") || die "can't find $programe\n";
6  my $program = '';
7  {
8      local $/ = undef;
9      $program = <$PROGRAM>;
10 }
11 $program =~ s/(?xs)
12     /\* # Match the opening delimiter
13     .*? # Match a minimal number of characters
14     \*/ # Match the closing delimiter
15     |
16     (?-s)//.* # C++ // comments. No s modifier
17 }[]g;
18
19 print $program;
```

Sigue un ejemplo de ejecución. Usaremos como entrada el programa C:

```
casiano@tonga:~/Lperltesting$ cat -n ehello.c
 1  #include <stdio.h>
 2  /* first
 3  comment
 4  */
 5  main() { // A C++ comment
 6      printf("hello world!\n"); /* second comment */
 7  } // final comment
```

Al ejecutar el programa eliminamos los comentarios:

```
casiano@tonga:~/Lperltesting$ extendedcomments.pl ehello.c | cat -n
 1  #include <stdio.h>
 2
 3  main() {
 4      printf("hello world!\n");
 5  }
```

3.2.3. Mirando hacia adetrás y hacia adelante

El siguiente fragmento esta 'casi' literalmente tomado de la sección 'Looking-ahead-and-looking-behind' en perlretut:

Las zero-width assertions como caso particular de mirar atrás-adelante

In Perl regular expressions, most regexp elements 'eat up' a certain amount of string when they match. For instance, the regexp element [abc] eats up one character of the string when it matches, in the sense that Perl moves to the next character position in the string after the match. There are some elements, however, that don't eat up characters (advance the character position) if they match.

The examples we have seen so far are the anchors. The anchor ^ matches the beginning of the line, but doesn't eat any characters.

Similarly, the word boundary anchor \b matches wherever a character matching \w is next to a character that doesn't, but it doesn't eat up any characters itself.

Anchors are examples of zero-width assertions. Zero-width, because they consume no characters, and assertions, because they test some property of the string.

In the context of our walk in the woods analogy to regexp matching, most regexp elements move us along a trail, but anchors have us stop a moment and check our surroundings. If the local environment checks out, we can proceed forward. But if the local environment doesn't satisfy us, we must backtrack.

Checking the environment entails either looking ahead on the trail, looking behind, or both.

- *^ looks behind, to see that there are no characters before.*
- *\$ looks ahead, to see that there are no characters after.*
- *\b looks both ahead and behind, to see if the characters on either side differ in their "word-ness".*

The lookahead and lookbehind assertions are generalizations of the anchor concept. Lookahead and lookbehind are zero-width assertions that let us specify which characters we want to test for.

Lookahead assertion

The lookahead assertion is denoted by (?=regexp) and the lookbehind assertion is denoted by (?<=fixed-regexp).

En español, operador de "trailing" o "mirar-adelante" positivo. Por ejemplo, /\w+(?=\t)/ solo casa una palabra si va seguida de un tabulador, pero el tabulador no formará parte de \$&. Ejemplo:

```
> cat -n lookahead.pl
1  #!/usr/bin/perl
2
3  $a = "bugs the rabbit";
4  $b = "bugs the frog";
5  if ($a =~ m{bugs(?: the cat| the rabbit)}i) { print "$a matches. \${&} = \${&}\n"; }
6  else { print "$a does not match\n"; }
7  if ($b =~ m{bugs(?: the cat| the rabbit)}i) { print "$b matches. \${&} = \${&}\n"; }
8  else { print "$b does not match\n"; }
```

Al ejecutar el programa obtenemos:

```
> lookahead.pl
bugs the rabbit matches. ${&} = bugs
bugs the frog does not match
>
```

Some examples using the debugger⁴:

```
DB<1>          #012345678901234567890
DB<2> $x = "I catch the housecat 'Tom-cat' with catnip"
DB<3> print "(${&} (".pos($x).")\n" if $x =~ /cat(?:\s)/g
(cat) (20)          # matches 'cat' in 'housecat'

DB<5> $x = "I catch the housecat 'Tom-cat' with catnip" # To reset pos
DB<6> x @catwords = ($x =~ /(?:\s)cat\w+/g)
0  'catch'
1  'catnip'

DB<7>          #012345678901234567890123456789
DB<8> $x = "I catch the housecat 'Tom-cat' with catnip"
DB<9> print "(${&} (".pos($x).")\n" if $x =~ /\bcat\b/g
(cat) (29) # matches 'cat' in 'Tom-cat'

DB<10> $x = "I catch the housecat 'Tom-cat' with catnip"
DB<11> x  $x =~ /(?:\s)cat(?:\s)/
empty array
DB<12> # doesn't match; no isolated 'cat' in middle of $x
```

A hard RegEx problem

Véase el nodo A hard RegEx problem en PerlMonks. Un monje solicita:

Hi Monks,

I wanna to match this issues:

- 1. The string length is between 3 and 10*
- 2. The string ONLY contains [0-9] or [a-z] or [A-Z], but*
- 3. The string must contain a number AND a letter at least*

Pls help me check. Thanks

Solución:

⁴catnip: La nepeta cataria, también llamada menta de los gatos, de la familia del tomillo y la lavanda. Su perfume desencadena un comportamiento en el animal, similar al del celo

```
casiano@millo:~$ perl -wde 0
main::(-e:1): 0
DB<1> x 'aaa2a1' =~ /\A(?=.*[a-z])(?=.*\d)\w{3,10}\z/i
0 1
DB<2> x 'aaaaaa' =~ /\A(?=.*[a-z])(?=.*\d)\w{3,10}\z/i
empty array
DB<3> x '1111111' =~ /\A(?=.*[a-z])(?=.*\d)\w{3,10}\z/i
empty array
DB<4> x '1111111bbbb' =~ /\A(?=.*[a-z])(?=.*\d)\w{3,10}\z/i
empty array
DB<5> x '111bbbb' =~ /\A(?=.*[a-z])(?=.*\d)\w{3,10}\z/i
0 1
```

Los paréntesis lookahead and lookbehind no capturan

Note that the parentheses in `(?=regex)` and `(?<=regex)` are non-capturing, since these are zero-width assertions.

Limitaciones del lookbehind

Lookahead `(?=regex)` can match arbitrary regexps, but lookbehind `(?<=fixed-regexp)` only works for regexps of fixed width, i.e., a fixed number of characters long.

Thus `(?<=(ab|bc))` is fine, but `(?<=(ab))` is not.*

Negación de los operadores de lookahead y lookbehind

The negated versions of the lookahead and lookbehind assertions are denoted by `(?!regex)` and `(?<!=fixed-regexp)` respectively. They evaluate true if the regexps do not match:

```
$x = "foobar";
$x =~ /foo(?!bar)/; # doesn't match, 'bar' follows 'foo'
$x =~ /foo(?!baz)/; # matches, 'baz' doesn't follow 'foo'
$x =~ /(?!\\s)foo/; # matches, there is no \\s before 'foo'
```

Ejemplo: split con lookahead y lookbehind

Here is an example where a string containing blank-separated words, numbers and single dashes is to be split into its components.

Using `/\\s+/` alone won't work, because spaces are not required between dashes, or a word or a dash. Additional places for a split are established by looking ahead and behind:

```
casiano@tonga:~$ perl5.10.1 -wdE 0
main::(-e:1): 0
DB<1> $str = "one two - --6-8"
DB<2> x @toks = split / \\s+ | (?<=\\S) (?=-) | (?<=-) (?=\\S)/x, $str
0 'one'
1 'two'
2 '-'
3 '-'
4 '-'
5 6
6 '-'
7 8
```

Look Around en perlre

El siguiente párrafo ha sido extraído la sección 'Look-Around-Assertions' en `perlre`. Usémoslo como texto de repaso:

Look-around assertions are zero width patterns which match a specific pattern without including it in `$&`. Positive assertions match when their subpattern matches, negative assertions match when their subpattern fails. Look-behind matches text up to the current match position, look-ahead matches text following the current match position.

- `(?=pattern)`

A zero-width positive look-ahead assertion. For example, `/\w+(?=\t)/` matches a word followed by a tab, without including the tab in `$&`.

- `(?!pattern)`

A zero-width negative look-ahead assertion. For example `/foo(?!bar)/` matches any occurrence of `foo` that isn't followed by `bar`.

Note however that look-ahead and look-behind are NOT the same thing. You cannot use this for look-behind.

If you are looking for a `bar` that isn't preceded by a `foo`, `/(?!foo)bar/` will not do what you want.

That's because the `(?!foo)` is just saying that the next thing cannot be `foo` –and it's not, it's a `bar`, so `foobar` will match.

You would have to do something like `/(?!foo)...\bar/` for that.

We say "like" because there's the case of your `bar` not having three characters before it.

You could cover that this way: `/(?: (?!foo) ... | ^.{0,2})bar/`. Sometimes it's still easier just to say:

```
if (/bar/ && $' !~ /foo$/)
```

For look-behind see below.

- `(?<=pattern)`

A zero-width positive look-behind assertion.

For example, `/(?<=\t)\w+/` matches a word that follows a tab, without including the tab in `$&`. Works only for fixed-width look-behind.

- `\K`

There is a special form of this construct, called `\K`, which causes the regex engine to 'keep' everything it had matched prior to the `\K` and not include it in `$&`. This effectively provides variable length look-behind. The use of `\K` inside of another look-around assertion is allowed, but the behaviour is currently not well defined.

For various reasons `\K` may be significantly more efficient than the equivalent `(?<=...)` construct, and it is especially useful in situations where you want to efficiently remove something following something else in a string. For instance

```
s/(foo)bar/$1/g;
```

can be rewritten as the much more efficient

```
s/foo\Kbar//g;
```

Sigue una sesión con el depurador que ilustra la semántica del operador:

```
casiano@millo:~$ perl5.10.1 -wdE 0
main::(-e:1): 0
```

```
DB<1> print "& = <$&> 1 = <$1>\n" if "alphabet" =~ /([^\aeiou][a-z][\aeiou])[a-z]/
& = <phab> 1 = <pha>
```

```
DB<2> print "& = <$&> 1 = <$1>\n" if "alphabet" =~ /\K([^\aeiou][a-z][\aeiou])[a-z]/
& = <phab> 1 = <pha>
```

```
DB<3> print "& = <$&> 1 = <$1>\n" if "alphabet" =~ /([^\aeiou]\K[a-z][\aeiou])[a-z]/
& = <hab> 1 = <pha>
```



```

DB<4> print "& = <$&> 1 = <$1>\n" if "alphabet" =~ /([^aeiou][a-z]\K[aeiou])[a-z]/
& = <ab> 1 = <pha>
DB<5> print "& = <$&> 1 = <$1>\n" if "alphabet" =~ /([^aeiou][a-z][aeiou])\K[a-z]/
& = <b> 1 = <pha>
DB<6> print "& = <$&> 1 = <$1>\n" if "alphabet" =~ /([^aeiou][a-z][aeiou])[a-z]\K/
& = <> 1 = <pha>
DB<7> @a = "alphabet" =~ /([aeiou]\K[^aeiou])/g; print "$&\n"
t
DB<8> x @a
0 'al'
1 'ab'
2 'et'

```

Otro ejemplo: eliminamos los blancos del final en una cadena:

```

DB<23> $x = ' cadena entre blancos '
DB<24> ($y = $x) =~ s/.*\b\K.*//g
DB<25> p "<$y>"
< cadena entre blancos>

```

■ (?<!pattern)

A zero-width negative look-behind assertion.

For example /(?!bar)foo/ matches any occurrence of foo that does not follow bar.

Works only for fixed-width look-behind.

Veamos un ejemplo de uso. Se quiere sustituir las extensiones `.something` por `.txt` en cadenas que contienen una ruta a un fichero:

```

casiano@millor:~$ perl5.10.1 -wdE 0
main::(-e:1): 0
DB<1> ($b = $a = 'abc/xyz.something') =~ s{\.[^.]*$}{.txt}
DB<2> p $b
abc/xyz.txt
DB<3> ($b = $a = 'abc/xyz.something') =~ s/\.\K[^.]*/txt/;
DB<4> p $b
abc/xyz.txt
DB<5> p $a
abc/xyz.something

```

Véase también:

- `Regexp::Keep` por Jeff Pinyan
- El nodo *positive look behind regex mystery* en PerlMonks

Operador de predicción negativo: Última ocurrencia

Escriba una expresión regular que encuentre la última aparición de la cadena `foo` en una cadena dada.

```

DB<6> x ($a = 'foo foo bar bar foo bar bar') =~ /foo(?!.*foo)/g; print pos($a)."\n"
19
DB<7> x ($a = 'foo foo bar bar foo bar bar') =~ s/foo(?!.*foo)/\U$&/
0 1
DB<8> x $a
0 'foo foo bar bar FOO bar bar'

```

Diferencias entre mirar adelante negativo y mirar adelante con clase negada

Aparentemente el operador “mirar-adelante” negativo es parecido a usar el operador “mirar-adelante” positivo con la negación de una clase.

<code>/regexp(?![abc])/</code>	<code>/regexp(?=[^abc])/</code>
--------------------------------	---------------------------------

Sin embargo existen al menos dos diferencias:

- Una negación de una clase debe casar algo para tener éxito. Un ‘mirar-adelante’ negativo tiene éxito si, en particular no logra casar con algo. Por ejemplo:

`\d+(?!\.)` casa con `$a = '452'`, mientras que `\d+(?=[^.])` lo hace, pero porque 452 es 45 seguido de un carácter que no es el punto:

```
> cat lookaheadneg.pl
#!/usr/bin/perl

$a = "452";
if ($a =~ m{\d+(?=[^.] )}i) { print "$a casa clase negada. \${&} = \${&}\n"; }
else { print "$a no casa\n"; }
if ($a =~ m{\d+(?!\. )}i) { print "$a casa predicción negativa. \${&} = \${&}\n"; }
else { print "$b no casa\n"; }
nereida:~/perl/src> lookaheadneg.pl
452 casa clase negada. \${&} = 45
452 casa predicción negativa. \${&} = 452
```

- Una clase negada casa un único carácter. Un ‘mirar-adelante’ negativo puede tener longitud arbitraria.

AND y AND NOT

Otros dos ejemplos:

- `^(?![A-Z]*$)[a-zA-Z]*$`

casa con líneas formadas por secuencias de letras tales que no todas son mayúsculas. (Obsérvese el uso de las anclas).

- `^(?=.*?esto)(?=.*?eso)`

casan con cualquier línea en la que aparezcan `esto` y `eso`. Ejemplo:

```
> cat estoyeso.pl
#!/usr/bin/perl

my $a = shift;

if ($a =~ m{^(?=.*?esto)(?=.*?eso)}i) { print "$a matches.\n"; }
else { print "$a does not match\n"; }

>estoyeso.pl 'hola eso y esto'
hola eso y esto matches.
> estoyeso.pl 'hola esto y eso'
hola esto y eso matches.
> estoyeso.pl 'hola aquello y eso'
```

```
hola aquello y eso does not match
> estoyeso.pl 'hola esto y aquello'
hola esto y aquello does not match
```

El ejemplo muestra que la interpretación es que cada operador mirar-adelante se interpreta siempre a partir de la posición actual de búsqueda. La expresión regular anterior es básicamente equivalente a `(/esto/ && /eso/)`.

- `(?!000)(\d\d\d)`

casa con cualquier cadena de tres dígitos que no sea la cadena 000.

Lookahead negativo versus lookbehind

Nótese que el “mirar-adelante” negativo no puede usarse fácilmente para imitar un “mirar-atrás”, esto es, que no se puede imitar la conducta de `(?<!foo)bar` mediante algo como `(/?!foo)bar`. Tenga en cuenta que:

- Lo que dice `(?!foo)` es que los tres caracteres que siguen no puede ser `foo`.
- Así, `foo` no pertenece a `(?!foo)bar/`, pero `foobar` pertenece a `(?!foo)bar/` porque `bar` es una cadena cuyos tres siguientes caracteres son `bar` y no son `foo`.
- Si quisieramos conseguir algo parecido a `(?<!foo)bar` usando un lookahead negativo tendríamos que escribir algo así como `(?!foo)...bar/` que casa con una cadena de tres caracteres que no sea `foo` seguida de `bar` (pero que tampoco es exactamente equivalente):

```
pl@nereida:~/Lperltesting$ cat -n foobar.pl
 1 use v5.10;
 2 use strict;
 3
 4 my $a = shift;
 5
 6 for my $r (q{(?<!foo)bar}, q{(?!foo)bar}, q{(?!foo)...bar}) {
 7     if ($a =~ /$r/) {
 8         say "$a casa con $r"
 9     }
10     else {
11         say "$a no casa con $r"
12     }
13 }
```

- Al ejecutar con diferentes entradas el programa anterior vemos que la solución `q{(?!foo)...bar}` se aproxima mas a `q{(?<!foo)bar}`:

```
pl@nereida:~/Lperltesting$ perl5.10.1 foobar.pl foobar
foobar no casa con (?<!foo)bar
foobar casa con (?!foo)bar
foobar no casa con (?!foo)...bar
```

```
pl@nereida:~/Lperltesting$ perl5.10.1 foobar.pl bar
bar casa con (?<!foo)bar
bar casa con (?!foo)bar
bar no casa con (?!foo)...bar
```

Ejercicio 3.2.1. *Explique porqué `bar` casa con `(?<!foo)bar` pero no con `(?!foo)...bar`. ¿Sabría encontrar una expresión regular mas apropiada usando lookahead negativo?*

- En realidad, posiblemente sea mas legible una solución como:

```
if (/bar/ and $' !~ /foo$/)
```

o aún mejor (véase 3.1.4):

```
if (/bar/p && ${^PREMATCH} =~ /foo$/)
```

El siguiente programa puede ser utilizado para ilustrar la equivalencia:

```
pl@nereida:~/Lperltesting$ cat -n foobarprematch.pl
1  use v5.10;
2  use strict;
3
4  $_ = shift;
5
6  if (/bar/p && ${^PREMATCH} =~ /foo$/) {
7      say "$_ no cumple ".q{/bar/p && ${^PREMATCH} =~ /foo$/};
8  }
9  else {
10     say "$_ cumple ".q{/bar/p && ${^PREMATCH} =~ /foo$/};
11 }
12 if (/(<?!foo)bar/) {
13     say "$_ casa con (<?!foo)bar"
14 }
15 else {
16     say "$_ no casa con (<?!foo)bar"
17 }
```

Siguen dos ejecuciones:

```
pl@nereida:~/Lperltesting$ perl5.10.1 foobarprematch.pl bar
bar cumple /bar/p && ${^PREMATCH} =~ /foo$/
bar casa con (<?!foo)bar
pl@nereida:~/Lperltesting$ perl5.10.1 foobarprematch.pl foobar
foobar no cumple /bar/p && ${^PREMATCH} =~ /foo$/
foobar no casa con (<?!foo)bar
```

Ejercicios

Ejercicio 3.2.2. ■ *Escriba una sustitución que reemplaze todas las apariciones de foo por foo, usando \K o lookbehind*

- *Escriba una sustitución que reemplaze todas las apariciones de lookahead por look-ahead usando lookaheads y lookbehinds*
- *Escriba una expresión regular que capture todo lo que hay entre las cadenas foo y bar siempre que no se incluya la palabra baz*
- *¿Cuál es la salida?*

```
DB<1> x 'abc' =~ /(?(.)(.)(.))a(b)/
```

- *Se quiere poner un espacio en blanco después de la aparición de cada coma:*

```
s/,/, /g;
```

pero se quiere que la sustitución no tenga lugar si la coma esta incrustada entre dos dígitos.

- Se quiere poner un espacio en blanco después de la aparición de cada coma:

```
s/,/, /g;
```

pero se quiere que la sustitución no tenga lugar si la coma esta incrustada entre dos dígitos. Además se pide que si hay ya un espacio después de la coma, no se duplique

- ¿Cuál es la salida?

```
pl@nereida:~/Lperltesting$ cat -n ABC123.pl
 1 use warnings;
 2 use strict;
 3
 4 my $c = 0;
 5 my @p = ('^(ABC)(?!123)', '^(\\D*)(?!123)',);
 6
 7 for my $r (@p) {
 8     for my $s (qw{ABC123 ABC445}) {
 9         $c++;
10         print "$c: '$s' =~ /$r/ : ";
11         <>;
12         if ($s =~ /$r/) {
13             print " YES ($1)\n";
14         }
15         else {
16             print " NO\n";
17         }
18     }
19 }
```

3.2.4. Definición de Nombres de Patrones

Perl 5.10 introduce la posibilidad de definir subpatrones en una sección del patrón.

Lo que dice perlretut sobre la definición de nombres de patrones

Citando la sección *Defining named patterns* en el documento la sección 'Defining-named-patterns' en perlretut para perl5.10:

Some regular expressions use identical subpatterns in several places. Starting with Perl 5.10, it is possible to define named subpatterns in a section of the pattern so that they can be called up by name anywhere in the pattern. This syntactic pattern for this definition group is "(? (DEFINE) (?<name>pattern) ...)" An insertion of a named pattern is written as (?&name).

Veamos un ejemplo que define el lenguaje de los números en punto flotante:

```
pl@nereida:~/Lperltesting$ cat -n definingnamedpatterns.pl
 1 #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w
 2 use v5.10;
 3
 4 my $regex = qr{
 5     ^ (?<num>
 6         (?&osg)[\t\ ]* (? : (?&int)(?&dec)? | (?&dec) )
```

```

7      )
8      (?: [eE]
9      (?<exp> (?&osg)(?&int)) )?
10     $
11     (? (DEFINE)
12         (?<osg>[+-]?)          # optional sign
13         (?<int>\d++)           # integer
14         (?<dec>\.(?&int))      # decimal fraction
15     )
16 }x;
17
18 my $input = <>;
19 chomp($input);
20 my @r;
21 if (@r = $input =~ $regex) {
22     my $exp = $+{exp} || '';
23     say "$input matches: (num => '$+{num}', exp => '$exp')";
24 }
25 else {
26     say "does not match";
27 }

```

perlretut comenta sobre este ejemplo:

The example above illustrates this feature. The three subpatterns that are used more than once are the optional sign, the digit sequence for an integer and the decimal fraction. The DEFINE group at the end of the pattern contains their definition. Notice that the decimal fraction pattern is the first place where we can reuse the integer pattern.

Lo que dice perlre sobre la definición de patrones

Curiosamente, (DEFINE) se considera un caso particular de las expresiones regulares condicionales de la forma (? (condition) yes-pattern) (véase la sección 3.2.10). Esto es lo que dice la sección 'Extended-Patterns' en perlre al respecto:

A special form is the (DEFINE) predicate, which never executes directly its yes-pattern, and does not allow a no-pattern. This allows to define subpatterns which will be executed only by using the recursion mechanism. This way, you can define a set of regular expression rules that can be bundled into any pattern you choose.

It is recommended that for this usage you put the DEFINE block at the end of the pattern, and that you name any subpatterns defined within it.

Also, it's worth noting that patterns defined this way probably will not be as efficient, as the optimiser is not very clever about handling them.

An example of how this might be used is as follows:

```

1. /(?(<NAME>(?(&NAME_PAT)))(?(<ADDR>(?(&ADDRESS_PAT)))
2.      (?(DEFINE)
3.          (?<NAME_PAT>....)
4.          (?<ADDRESS_PAT>....)
5. )/x

```

Note that capture buffers matched inside of recursion are not accessible after the recursion returns, so the extra layer of capturing buffers is necessary. Thus \$+{NAME_PAT} would not be defined even though \$+{NAME} would be.

Lo que dice perlvar sobre patrones con nombre Esto es lo que dice perlvar respecto a las variables implicadas %+ y %-. Con respecto a el hash %+:

■ **%LAST_PAREN_MATCH, %+**

Similar to @+ , the %+ hash allows access to the named capture buffers, should they exist, in the last successful match in the currently active dynamic scope.

For example, \${foo} is equivalent to \$1 after the following match:

```
1. 'foo' =~ /(?(<foo>foo)/;
```

The keys of the %+ hash list only the names of buffers that have captured (and that are thus associated to defined values).

The underlying behaviour of %+ is provided by the Tie::Hash::NamedCapture module.

Note: %- and %+ are tied views into a common internal hash associated with the last successful regular expression. Therefore mixing iterative access to them via each may have unpredictable results. Likewise, if the last successful match changes, then the results may be surprising.

■ **%-**

*Similar to %+ , this variable allows access to the named capture buffers in the last successful match in the currently active dynamic scope. **To each capture buffer name found in the regular expression, it associates a reference to an array containing the list of values captured by all buffers with that name (should there be several of them), in the order where they appear.***

Here's an example:

```
1. if ('1234' =~ /(?(<A>1)(?(B>2)(?(A>3)(?(B>4)/) {  
2.   foreach my $bufname (sort keys %-) {  
3.     my $ary = ${$bufname};  
4.     foreach my $idx (0..$#$ary) {  
5.       print "\${$bufname}[$idx] : ",  
6.         (defined($ary->[$idx]) ? "'$ary->[$idx]'" : "undef"),  
7.         "\n";  
8.     }  
9.   }  
10. }
```

would print out:

```
1. ${A}[0] : '1'  
2. ${A}[1] : '3'  
3. ${B}[0] : '2'  
4. ${B}[1] : '4'
```

The keys of the %- hash correspond to all buffer names found in the regular expression.

3.2.5. Patrones Recursivos

Perl 5.10 introduce la posibilidad de definir subpatrones en una sección del patrón. Citando la versión del documento perlretut para perl5.10:

This feature (introduced in Perl 5.10) significantly extends the power of Perl's pattern matching. By referring to some other capture group anywhere in the pattern with the construct (?group-ref), the pattern within the referenced group is used as an independent subpattern in place of the group reference itself. Because the group reference may be contained within the group it refers to, it is now possible to apply pattern matching to tasks that hitherto required a recursive parser.

...
In (?...) both absolute and relative backreferences may be used. The entire pattern can be reinserted with (?R) or (?0). If you prefer to name your buffers, you can use (?&name) to recurse into that buffer.

Palíndromos

Véase un ejemplo que reconoce los palabra-palíndromos (esto es, la lectura directa y la inversa de la cadena pueden diferir en los signos de puntuación):

```
casiano@millo:~/Lperltesting$ cat -n palindromos.pl
 1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w
 2  use v5.10;
 3
 4  my $regex = qr/^(\\W*
 5              (?:
 6                  (\\w) (?1) \\g{-1}  # palindromo estricto
 7                  |
 8                  \\w?                # no recursiva
 9              )
10              \\W*)$/ix;
11
12  my $input = <>;
13  chomp($input);
14  if ($input =~ $regex) {
15      say "$input is a palindrome";
16  }
17  else {
18      say "does not match";
19  }
```

Ejercicio 3.2.3. *¿Cuál es el efecto del modificador i en la regexp `qr/^(\\W* (?: (\\w) (?1) \\g{-1} | \\w?) \\W*`*

Siguen algunos ejemplos de ejecución⁵

```
pl@nereida:~/Lperltesting$ ./palindromos.pl
A man, a plan, a canal: Panama!
A man, a plan, a canal: Panama! is a palindrome
pl@nereida:~/Lperltesting$ ./palindromos.pl
A man, a plan, a cam, a yak, a yam, a canal { Panama!
A man, a plan, a cam, a yak, a yam, a canal { Panama! is a palindrome
pl@nereida:~/Lperltesting$ ./palindromos.pl
A man, a plan, a cat, a ham, a yak, a yam, a hat, a canal { Panama!
A man, a plan, a cat, a ham, a yak, a yam, a hat, a canal { Panama! is a palindrome
pl@nereida:~/Lperltesting$ ./palindromos.pl
saippuakauppias
saippuakauppias is a palindrome
pl@nereida:~/Lperltesting$ ./palindromos.pl
dfghjgfd
does not match
```

⁵

- saippuakauppias: Vendedor de jabón (suomi)
- yam: batata (inglés)
- cam: leva


```
pl@nereida:~/Lperltesting$ ./palindromos.pl
...;...;
...;...; is a palindrome
```

Lo que dice perlre sobre recursividad

```
(?PARNO) (?-PARNO) (?+PARNO) (?R) (?0)
```

Similar to (??{ code }) (véase la sección 3.2.9) except it does not involve compiling any code, instead it treats the contents of a capture buffer as an independent pattern that must match at the current position. Capture buffers contained by the pattern will have the value as determined by the outermost recursion.

PARNO is a sequence of digits (not starting with 0) whose value reflects the paren-number of the capture buffer to recurse to.

(?R) recurses to the beginning of the whole pattern. (?0) is an alternate syntax for (?R).

If PARNO is preceded by a plus or minus sign then it is assumed to be relative, with negative numbers indicating preceding capture buffers and positive ones following. Thus (?-1) refers to the most recently declared buffer, and (?+1) indicates the next buffer to be declared.

Note that the counting for relative recursion differs from that of relative backreferences, in that with recursion unclosed buffers are included.

Hay una diferencia fundamental entre `\g{-1}` y `(?-1)`. El primero significa *lo que pasó con el último paréntesis*. El segundo significa que se debe llamar a la expresión regular que define el último paréntesis. Véase un ejemplo:

```
pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> x ($a = "12 aAbB 34") =~ s/([aA])(?-1)(?+1)([bB])/-\1\2-/g
0 1
DB<2> p $a
12 -aB- 34
```

En `perlre` también se comenta sobre este punto:

If there is no corresponding capture buffer defined, then it is a fatal error. Recursing deeper than 50 times without consuming any input string will also result in a fatal error. The maximum depth is compiled into perl, so changing it requires a custom build.

Paréntesis Equilibrados

El siguiente programa (inspirado en uno que aparece en `perlre`) reconoce una llamada a una función `foo()` que puede contener una secuencia de expresiones con paréntesis equilibrados como argumento:

```
1 pl@nereida:~/Lperltesting$ cat perlrebalancedpar.pl
2 #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w
3 use v5.10;
4 use strict;
5
6 my $regex = qr{ (                               # paren group 1 (full function)
7     foo
8     (                                           # paren group 2 (parens)
9         \ (
10             (                                   # paren group 3 (contents of parens)
11                 (?:
```

```

12             [^()]+ # Non-parens
13             |
14             (?2) # Recurse to start of paren group 2
15             )*
16         ) # 3
17     \)
18     ) # 2
19     ) # 1
20 }x;
21
22 my $input = <>;
23 chomp($input);
24 my @res = ($input =~ /$regexp/);
25 if (@res) {
26     say "<$&> is balanced\nParen: (@res)";
27 }
28 else {
29     say "does not match";
30 }

```

Al ejecutar obtenemos:

```

pl@nereida:~/Lperltesting$ ./perlrebalancedpar.pl
foo(bar(baz)+baz(bop))
<foo(bar(baz)+baz(bop))> is balanced
Paren: (foo(bar(baz)+baz(bop)) (bar(baz)+baz(bop)) bar(baz)+baz(bop))

```

Como se comenta en `perlre` es conveniente usar índices relativos si se quiere tener una expresión regular reciclable:

The following shows how using negative indexing can make it easier to embed recursive patterns inside of a `qr//` construct for later use:

```

1. my $parens = qr/(\((?:[^\)]++|(?-1))*\))/;
2. if (/foo $parens \s+ + \s+ bar $parens/x) {
3.     # do something here...
4. }

```

Véase la sección 3.2.6 para comprender el uso de los operadores posesivos como `++`.

Capturando los bloques de un programa

El siguiente programa presenta una heurística para determinar los bloques de un programa:

```

1  pl@nereida:~/Lperltesting$ cat blocks.pl
2  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w
3  use v5.10;
4  use strict;
5  #use re 'debug';
6
7  my $rb = qr{(?x)
8      (
9          \{ # llave abrir
10         (?:
11             [^{}]* # no llaves

```

```

12         |
13         [^{}]*+ # no llaves
14         (?1)    # recursivo
15         [^{}]*+ # no llaves
16     )**+
17     \}          # llave cerrar
18 )
19 };
20
21 local $/ = undef;
22 my $input = <>;
23 my@blocks = $input =~ m{$rb}g;
24 my $i = 0;
25 say($i++."\n$_\n===") for @blocks;

```

Veamos una ejecución. Le daremos como entrada el siguiente programa: Al ejecutar el programa con esta entrada obtenemos:

<pre> pl@nereida:~/Lperltesting\$ cat -n blocks.pl 1 main() { /* 1 */ 2 { /* 2 */ } 3 { /* 3 */ } 4 } 5 6 f(){ /* 4 */ 7 { /* 5 */ 8 { /* 6 */ } 9 } 10 { /* 7 */ 11 { /* 8 */ } 12 } 13 } 14 15 g(){ /* 9 */ 16 } 17 18 h() { 19 {{{}}} 20 } 21 /* end h */ </pre>	<pre> pl@nereida:~/Lperltesting\$ perl5.10.1 blocks.pl 0: { /* 1 */ { /* 2 */ } { /* 3 */ } } === 1: { /* 4 */ { /* 5 */ { /* 6 */ } } { /* 7 */ { /* 8 */ } } } === 2: { /* 9 */ } === 3: { {{{}}} } === </pre>
---	--

Reconocimiento de Lenguajes Recursivos: Un subconjunto de L^AT_EX

La posibilidad de combinar en las expresiones regulares Perl 5.10 la recursividad con los constructos (`?<name>...`) y `?&name`) así como las secciones (`? (DEFINE) ...`) permiten la escritura de expresiones regulares que reconocen lenguajes recursivos. El siguiente ejemplo muestra un reconocedor de un subconjunto del lenguaje L^AT_EX (véase la entrada LaTeX en la wikipedia):

```

1 pl@nereida:~/Lperltesting$ cat latex5_10.pl
2 #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w

```

```

3 use strict;
4 use v5.10;
5
6 my $regexp = qr{
7     \A(?&File)\z
8
9     (? (DEFINE)
10         (?<File>      (?&Element)**\s*
11             )
12
13         (?<Element>  \s* (?&Command)
14             | \s* (?&Literal)
15         )
16
17         (?<Command>  \\ \s* (?<L>(?&Literal)) \s* (?<Op>(?&Options)?) \s* (?<A>(?&Args))
18             (?{
19                 say "command: <${L}> options: <${Op}> args: <${A}>"
20             })
21         )
22
23         (?<Options>  \[ \s* (?: (?&Option) (?: \s*, \s* (?&Option) )*)? \s* \]
24         )
25
26         (?<Args>     (?: \{ \s* (?&Element)* \s* \} ) *
27         )
28
29         (?<Option>   \s* [^ [\${&%#_}{~}\s,]+
30         )
31
32         (?<Literal>  \s* ([^ [\${&%#_}{~}\s,]+)
33         )
34     )
35 }xms;
36
37 my $input = do{ local $/; <> };
38 if ($input =~ $regexp) {
39     say "$@: matches:\n$&";
40 }
41 else {
42     say "does not match";
43 }

```

Añadimos una acción semántica al final de la aceptación de un <Command>.

```

    (?<Command>  \\ \s* (?<L>(?&Literal)) \s* (?<Op>(?&Options)?) \s* (?<A>(?&Args)?)
    (?{
        say "command: <${L}> options: <${Op}> args: <${A}>"
    })
)

```

Esta acción es ejecutada pero no afecta al proceso de análisis. (véase la sección 3.2.8 para mas información sobre las acciones semánticas en medio de una regexp). La acción se limita a mostrar que ha casado con cada una de las tres componentes: el comando, las opciones y los argumentos.

Los paréntesis adicionales, como en (?<L>(?&Literal)) son necesarios para guardar lo que casó.

Cuando se ejecuta produce la siguiente salida⁶:

```
pl@nereida:~/Lperltesting$ cat prueba.tex
\documentclass[a4paper,11pt]{article}
\usepackage{latexsym}
\author{D. Conway}
\title{Parsing \LaTeX{}}
\begin{document}
\maketitle
\tableofcontents
\section{Description}
...is easy \footnote{But not\\ \emph{necessarily} simple}.
In fact it's easy peasy to do.
\end{document}

pl@nereida:~/Lperltesting$ ./latex5_10.pl prueba.tex
command: <documentclass> options: <[a4paper,11pt]> args: <{article}>
command: <usepackage> options: <> args: <{latexsym}>
command: <author> options: <> args: <{D. Conway}>
command: <LaTeX> options: <> args: <{}>
command: <title> options: <> args: <{Parsing \LaTeX{}}>
command: <begin> options: <> args: <{document}>
command: <maketitle> options: <> args: <>
command: <tableofcontents> options: <> args: <>
command: <section> options: <> args: <{Description}>
command: <emph> options: <> args: <{necessarily}>
command: <footnote> options: <> args: <{But not\\ \emph{necessarily} simple}>
command: <end> options: <> args: <{document}>
: matches:
\documentclass[a4paper,11pt]{article}
\usepackage{latexsym}
\author{D. Conway}
\title{Parsing \LaTeX{}}
\begin{document}
\maketitle
\tableofcontents
\section{Description}
...is easy \footnote{But not\\ \emph{necessarily} simple}.
In fact it's easy peasy to do.
```

6

- peasy: A disagreeable taste of very fresh green peas
- easy peasy:
 1. (uk) very easy (short for easy-peasy-lemon-squeezy)
 2. the first half of a rhyming phrase with several alternate second halves, all of which connote an activity or a result that is, respectively, simple to perform or achieve.

Tie your shoes? Why that's easy peasy lemon squeezy!
Beat your meat? Why that's easy peasy Japanesey!
As a red-stater, condemn books and films without having read or seen them? Why that's easy peasy puddin'n'pie!
 3. It comes from a 1970's british TV commercial for Lemon Squeezy detergent. They were with a little girl who points out dirty greasy dishes to an adult (mom or relative) and then this adult produces Lemon Squeezy and they clean the dishes quickly. At the end of the commercial the girl says *Easy Peasy Lemon Squeezy*. Today it is a silly way to state something was or will be very easy.

```
\end{document}
```

La siguiente entrada `prueba3.tex` no pertenece al lenguaje definido por el patrón regular, debido a la presencia de la cadena `In` en la última línea:

```
pl@nereida:~/Lperltesting$ cat prueba3.tex
\documentclass[a4paper,11pt]{article}
\usepackage{latexsym}
\author{D. Conway}
\title{Parsing \LaTeX{}}
\begin{document}
\maketitle
\tableofcontents
\section{Description}
\comm{a}{b}
...is easy \footnote{But not\\ \emph{necessarily} simple}.
$In$ fact it's easy peasy to do.
\end{document}
```

```
pl@nereida:~/Lperltesting$ ./latex5_10.pl prueba3.tex
command: <documentclass> options: <[a4paper,11pt]> args: <{article}>
command: <usepackage> options: <> args: <{latexsym}>
command: <author> options: <> args: <{D. Conway}>
command: <LaTeX> options: <> args: <{}>
command: <title> options: <> args: <{Parsing \LaTeX{}}>
command: <begin> options: <> args: <{document}>
command: <maketitle> options: <> args: <>
command: <tableofcontents> options: <> args: <>
command: <section> options: <> args: <{Description}>
command: <comm> options: <> args: <{a}{b}>
command: <emph> options: <> args: <{necessarily}>
command: <footnote> options: <> args: <{But not\\ \emph{necessarily} simple}>
does not match
```

Ejercicio 3.2.4. *Obsérvese el uso del cuantificador posesivo en:*

```
10      (?<File>      (?&Element))*+\s*
11      )
```

¿Que ocurre si se quita el posesivo y se vuelve a ejecutar `$./latex5_10.pl prueba3.tex`?

Reconocimiento de Expresiones Aritméticas

Véase el nodo `Complex regex for maths formulas` en `perlmonks` para la formulación del problema. Un monje pregunta:

Hiya monks,

Im having trouble getting my head around a regular expression to match sequences. I need to catch all exceptions where a mathematical expression is illegal...

There must be either a letter or a digit either side of an operator parenthesis must open and close next to letters or digits, not next to operators, and do not have to exist variables must not be more than one letter Nothing other than a-z,A-Z,0-9,+,-,,/, (,) can be used*

*Can anyone offer a hand on how best to tackle this problem?
many thanks*

La solución parte de que una *expresión* es o bien un *término* o bien un *término* seguido de una operador y un *término*, esto es:

- termino
- termino op termino op termino ...

que puede ser unificado como `termino (op termino)*`.

Un *término* es un número o un identificador o una *expresión* entre paréntesis, esto es:

- numero
- identificador
- (expresión)

La siguiente expresión regular recursiva sigue esta idea:

```
pl@nereida:~/Lperltesting$ cat -n simpleexpressionsna.pl
 1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
 2  use v5.10;
 3  use strict;
 4  use warnings;
 5
 6  local our ($skip, $term, $expr);
 7  $skip = qr/\s*/;
 8  $expr = qr{ (?<EXPR>
 9              (?<TERM>          # An expression is a TERM ...
10                  $skip (?<ID>[a-zA-Z]+)
11                  | $skip (?<INT>[1-9]\d*)
12                  | $skip \(
13                      $skip (?&EXPR)
14                      $skip \)
15                  ) (? : $skip          # possibly followed by a sequence of ...
16                      (?<OP>[-+*/])
17                      (?&TERM)          # ... operand TERM pairs
18                  )*
19              )
20      }x;
21  my $re = qr/^ $expr $skip \z/x;
22  sub is_valid { shift =~ /$re/o }
23
24  my @test = ( '(a + 3)', '(3 * 4)+(b + x)', '(5 - a)*z',
25              '((5 - a))*(((z))) + 2)', '3 + 2', '!3 + 2', '3 + 2!',
26              '3 a', '3 3', '3 * * 3',
27              '2 - 3 * 4', '2 - 3 + 4',
28              );
29  foreach (@test) {
30      say("$_:");
31      say(is_valid($_) ? "\n<$_> is valid" : "\n<$_> is not valid")
32  }
```

Podemos usar acciones semánticas empotradas para ver la forma en la que trabaja la expresión regular (véase la sección 3.2.8):

```

pl@nereida:~/Lperltesting$ cat -n simpleexpressions.pl
 1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
 2  use v5.10;
 3  use strict;
 4  use warnings;
 5
 6  use re 'eval'; # to allow Eval-group at runtime
 7
 8  local our ($skip, $term, $expr);
 9  $skip = qr/\s*/;
10  $expr = qr{ (?<EXPR>
11              (?<TERM>                                # An expression is a TERM ...
12                  $skip (?<ID>[a-zA-Z]+) (?{ print "[ID ${ID}] " })
13                  | $skip (?<INT>[1-9]\d*) (?{ print "[INT ${INT}] " })
14                  | $skip \(                          (?{ print "([ " })
15                  $skip (?&EXPR)
16                  $skip \)                             (?{ print "[)] " })
17              ) (? : $skip                             # possibly followed by a sequence of ...
18                  (?<OP>[-+*/])                         (?{ print "[OP ${OP}] " })
19                  (?&TERM)                               # ... operand TERM pairs
20              )*
21          )
22      }x;
23  my $re = qr/^ $expr $skip \z/x;
24  sub is_valid { shift =~ /$re/o }
25
26  my @test = ( '(a + 3)', '(3 * 4)+(b + x)', '(5 - a)*z',
27              '((5 - a))*(((z))) + 2)', '3 + 2', '!3 + 2', '3 + 2!',
28              '3 a', '3 3', '3 * * 3',
29              '2 - 3 * 4', '2 - 3 + 4',
30          );
31  foreach (@test) {
32      say("$_:");
33      say(is_valid($_) ? "\n<$_> is valid" : "\n<$_> is not valid")
34  }

```

Ejecución:

```

pl@nereida:~/Lperltesting$ ./simpleexpressions.pl
(a + 3):
[()] [ID a] [OP +] [INT 3] []]
<(a + 3)> is valid
(3 * 4)+(b + x):
[()] [INT 3] [OP *] [INT 4] []] [OP +] [()] [ID b] [OP +] [ID x] []]
<(3 * 4)+(b + x)> is valid
(5 - a)*z:
[()] [INT 5] [OP -] [ID a] []] [OP *] [ID z]
<(5 - a)*z> is valid
((5 - a))*(((z))) + 2):
[()] [()] [INT 5] [OP -] [ID a] []] []] [OP *] [()] [()] [()] [ID z] []] []] [OP +] [INT 2]
<((5 - a))*(((z))) + 2)> is valid
3 + 2:
[INT 3] [OP +] [INT 2]
<3 + 2> is valid

```



```

!3 + 2:

<!3 + 2> is not valid
3 + 2!:
[INT 3] [OP +] [INT 2]
<3 + 2!> is not valid
3 a:
[INT 3]
<3 a> is not valid
3 3:
[INT 3]
<3 3> is not valid
3 * * 3:
[INT 3] [OP *]
<3 * * 3> is not valid
2 - 3 * 4:
[INT 2] [OP -] [INT 3] [OP *] [INT 4]
<2 - 3 * 4> is valid
2 - 3 + 4:
[INT 2] [OP -] [INT 3] [OP +] [INT 4]
<2 - 3 + 4> is valid

```

3.2.6. Cuantificadores Posesivos

Por defecto, cuando un subpatrón con un cuantificador impide que el patrón global tenga éxito, se produce un backtrack. Hay ocasiones en las que esta conducta da lugar a ineficiencia.

Perl 5.10 provee los cuantificadores posesivos: Un cuantificador posesivo actúa como un cuantificador greedy pero no se produce backtracking.

++	Casar 0 o mas veces y no retroceder
++	Casar 1 o mas veces y no retroceder
?+	Casar 0 o 1 veces y no retroceder
{n}+	Casar exactamente n veces y no retroceder (redundante)
{n,}+	Casar al menos n veces y no retroceder
{n,m}+	Casar al menos n veces y no mas de m veces y no retroceder

Por ejemplo, la ca-

dena 'aaaa' no casa con /(a++a)/ porque no hay retroceso después de leer las 4 aes:

```

pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> x 'aaaa' =~ /(a+a)/
0 'aaaa'
DB<2> x 'aaaa' =~ /(a++a)/
empty array

```

Cadenas Delimitadas por Comillas Dobles

Los operadores posesivos sirven para poder escribir expresiones regulares mas eficientes en aquellos casos en los que sabemos que el retroceso no conducirá a nuevas soluciones, como es el caso del reconocimiento de las cadenas delimitadas por comillas dobles:

```

pl@nereida:~/Lperltesting$ cat -n ./quotedstrings.pl
 1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
 2  use v5.10;
 3
 4  my $regex = qr/
 5      "          # double quote

```

```

6      (? :          # no memory
7          [^"\\]++  # no " or escape: Don't backtrack
8          | \\.      # escaped character
9      )*+
10     "              # end double quote
11    /x;
12
13    my $input = <>;
14    chomp($input);
15    if ($input =~ $regex) {
16        say "$& is a string";
17    }
18    else {
19        say "does not match";
20    }

```

Paréntesis Posesivos

Los paréntesis posesivos (`?> ...`) dan lugar a un reconocedor que rechaza las demandas de retroceso. De hecho, los operadores posesivos pueden ser reescritos en términos de los paréntesis posesivos: La notación `X++` es equivalente a `(?>X+)`.

Paréntesis Balanceados

El siguiente ejemplo reconoce el lenguaje de los paréntesis balanceados:

```

pl@nereida:~/Lperltesting$ cat -n ./balancedparenthesis.pl
1  #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1
2  use v5.10;
3
4  my $regex =
5      qr/^(
6          [^()]*+ # no hay paréntesis, no backtrack
7          \(
8              (?>      # subgrupo posesivo
9                  [^()]*+ # no hay paréntesis, + posesivo, no backtrack
10                 | (?1)  # o es un paréntesis equilibrado
11             )*
12          \)
13          [^()]*+ # no hay paréntesis
14      )$/x;
15
16  my $input = <>;
17  chomp($input);
18  if ($input =~ $regex) {
19      say "$& is a balanced parenthesis";
20  }
21  else {
22      say "does not match";
23  }

```

Cuando se ejecuta produce una salida como:

```

pl@nereida:~/Lperltesting$ ./balancedparenthesis.pl
(2*(3+4)-5)*2
(2*(3+4)-5)*2 is a balanced parenthesis

```

```

pl@nereida:~/Lperltesting$ ./balancedparenthesis.pl
(2*(3+4)-5))*2
does not match
pl@nereida:~/Lperltesting$ ./balancedparenthesis.pl
2*(3+4
does not match
pl@nereida:~/Lperltesting$ ./balancedparenthesis.pl
4*(2*(3+4)-5)*2
4*(2*(3+4)-5)*2 is a balanced parenthesis

```

Encontrando los bloques de un programa

El uso de los operadores posesivos nos permite reescribir la solución al problema de encontrar los bloques maximales de un código dada en la sección 3.2.5 de la siguiente manera:

```

1 pl@nereida:~/Lperltesting$ cat blocksopti.pl
2 #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w
3 use v5.10;
4 use strict;
5 #use re 'debug';
6
7 my $rb = qr{(?x)
8     (
9         \{                # llave abrir
10            (?
11                [^{}]+    # no llaves
12                |
13                (?1)      # recursivo
14                [^{}]*+    # no llaves
15            )*+
16        \}                # llave cerrar
17    )
18 };
19
20 local $/ = undef;
21 my $input = <>;
22 my@blocks = $input =~ m{$rb}g;
23 my $i = 0;
24 say($i++.":\n$_\n===") for @blocks;

```

Véase también

- Possessive Quantifiers en <http://www.regular-expressions.info/>
- Nodo *Possessive Quantifiers in Perl 5.10 regexps* en PerlMonks
- `perldoc perlre`

3.2.7. Perl 5.10: Numeración de los Grupos en Alternativas

A veces conviene tener una forma de acceso uniforme a la lista proporcionada por los paréntesis con memoria. Por ejemplo, la siguiente expresión regular reconoce el lenguaje de las horas en notaciones civil y militar:

```

pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0

```

```
DB<1> '23:12' =~ /(\d\d|\d):(\d\d)|(\d\d)(\d\d)/; print "1->$1 2->$2\n"
1->23 2->12
```

```
DB<2> '2312' =~ /(\d\d|\d):(\d\d)|(\d\d)(\d\d)/; print "3->$3 4->$4\n"
3->23 4->12
```

Parece inconveniente tener los resultados en variables distintas. El constructo `(?| ...)` hace que los paréntesis se enumeren relativos a las alternativas:

```
DB<3> '2312' =~ /(?(?|(\d\d|\d):(\d\d)|(\d\d)(\d\d)))/; print "1->$1 2->$2\n"
1->23 2->12
```

```
DB<4> '23:12' =~ /(?(?|(\d\d|\d):(\d\d)|(\d\d)(\d\d)))/; print "1->$1 2->$2\n"
1->23 2->12
```

Ahora en ambos casos \$1 y \$2 contienen las horas y minutos.

3.2.8. Ejecución de Código dentro de una Expresión Regular

Es posible introducir código Perl dentro de una expresión regular. Para ello se usa la notación `(?{code})`.

El siguiente texto esta tomado de la sección 'A-bit-of-magic:-executing-Perl-code-in-a-regular-expression' en `perlretut`:

Normally, regexps are a part of Perl expressions. Code evaluation expressions turn that around by allowing arbitrary Perl code to be a part of a regexp. A code evaluation expression is denoted `(?code)`, with code a string of Perl statements.

Be warned that this feature is considered experimental, and may be changed without notice.

Code expressions are zero-width assertions, and the value they return depends on their environment.

There are two possibilities: either the code expression is used as a conditional in a conditional expression `(?(condition)...)...`, or it is not.

- *If the code expression is a conditional, the code is evaluated and the result (i.e., the result of the last statement) is used to determine truth or falsehood.*
- *If the code expression is not used as a conditional, the assertion always evaluates true and the result is put into the special variable `$^R`. The variable `$^R` can then be used in code expressions later in the regexp*

Resultado de la última ejecución

Las expresiones de código son *zero-width assertions*: no consumen entrada. El resultado de la ejecución se salva en la variable especial `$^R`.

Veamos un ejemplo:

```
pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> $x = "abcdef"
DB<2> $x =~ /abc(?{ "Hi mom\n" })def(?{ print $^R })/
Hi mom
DB<3> $x =~ /abc(?{ print "Hi mom\n"; 4 })def(?{ print "$^R\n" })/
Hi mom
4
DB<4> $x =~ /abc(?{ print "Hi mom\n"; 4 })ddd(?{ print "$^R\n" })/ # does not match
DB<5>
```

En el último ejemplo (línea DB<4>) ninguno de los `print` se ejecuta dado que no hay matching.

El Código empotrado no es interpolado

Tomado de la sección 'Extended-Patterns' en `perlre`:

This zero-width assertion evaluates any embedded Perl code. It always succeeds, and its code is not interpolated. Currently, the rules to determine where the code ends are somewhat convoluted.

Contenido del último paréntesis y la variable por defecto en acciones empotradas

Tomado de la sección 'Extended-Patterns' en `perlre`:

... can be used with the special variable `$^N` to capture the results of submatches in variables without having to keep track of the number of nested parentheses. For example:

```
pl@nereida:~/Lperltesting$ perl5.10.1 -wdE 0
main::(-e:1): 0
DB<1> $x = "The brown fox jumps over the lazy dog"
DB<2> x $x =~ /the (\S+)(?{ $color = $^N }) (\S+)(?{ $animal = $^N })/i
0 'brown'
1 'fox'
DB<3> p "color=$color animal=$animal\n"
color=brown animal=fox
DB<4> $x =~ /the (\S+)(?{ print (substr($_,0,pos($_)))."\n" }) (\S+)/i
The brown
```

Inside the `(?{...})` block, `$_` refers to the string the regular expression is matching against. You can also use `pos()` to know what is the current position of matching within this string.

Los cuantificadores y el código empotrado

Si se usa un cuantificador sobre un código empotrado, actúa como un bucle:

```
pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> $x = "aaaa"
DB<2> $x =~ /(a(?{ $c++ }))*/
DB<3> p $c
4
DB<4> $y = "abcd"
DB<5> $y =~ /(?:.)(?{ print "-$1-\n" }))*/
-a-
-b-
-c-
-d-
```

Ámbito

Tomado (y modificado el ejemplo) de la sección 'Extended-Patterns' en `perlre`:

... The code is properly scoped in the following sense: If the assertion is backtracked (compare la sección 'Backtracking' en `perlre`), all changes introduced after localization are undone, so that

```
pl@nereida:~/Lperltesting$ cat embededcodescope.pl
use strict;

our ($cnt, $res);
```

```

sub echo {
local our $pre = substr($_,0,pos($_));
local our $post = (pos($_) < length)? (substr($_,1+pos($_))) : '';

print("$pre(count = $cnt)$post\n");
}

$_ = 'a' x 8;
m<
(?{ $cnt = 0 }) # Initialize $cnt.
(
  a
  (?{
    local $cnt = $cnt + 1; # Update $cnt, backtracking-safe.
    echo();
  })
)*
aaaa
(?{ $res = $cnt }) # On success copy to non-localized
# location.
>x;

print "FINAL RESULT: cnt = $cnt res =$res\n";

```

will set \$res = 4 . Note that after the match, \$cnt returns to the globally introduced value, because the scopes that restrict local operators are unwound.

```

pl@nereida:~/Lperltesting$ perl5.8.8 -w embedcodescope.pl
a(count = 1)aaaaaa
aa(count = 2)aaaaa
aaa(count = 3)aaaa
aaaa(count = 4)aaa
aaaaa(count = 5)aa
aaaaaa(count = 6)a
aaaaaaa(count = 7)
aaaaaaaa(count = 8)
FINAL RESULT: cnt = 0 res =4

```

Caveats

- *Due to an unfortunate implementation issue, the Perl code contained in these blocks is treated as a compile time closure that can have seemingly bizarre consequences when used with lexically scoped variables inside of subroutines or loops. There are various workarounds for this, including simply using global variables instead. If you are using this construct and strange results occur then check for the use of lexically scoped variables.*
- *For reasons of security, this construct is forbidden if the regular expression involves run-time interpolation of variables, unless the perilous `use re 'eval'` pragma has been used (see `re`), or the variables contain results of `qr//` operator (see "`qr/STRING/imosx`" in `perlop`).*
This restriction is due to the wide-spread and remarkably convenient custom of using run-time determined strings as patterns. For example:

1. `$re = <>;`
2. `chomp $re;`
3. `$string =~ /$re/;`

Before Perl knew how to execute interpolated code within a pattern, this operation was completely safe from a security point of view, although it could raise an exception from an illegal pattern. If you turn on the `use re 'eval'`, though, it is no longer secure, so you should only do so if you are also using taint checking. Better yet, use the carefully constrained evaluation within a Safe compartment. See `perlsec` for details about both these mechanisms. (Véase la sección 'Taint-mode' en `perlsec`)

- *Because Perl's regex engine is currently not re-entrant, interpolated code may not invoke the regex engine either directly with `m//` or `s///`, or indirectly with functions such as `split`.*

Depurando con código empotrado Colisiones en los Nombres de las Subexpresiones Regulares

Las acciones empotradas pueden utilizarse como mecanismo de depuración y de descubrimiento del comportamiento de nuestras expresiones regulares.

En el siguiente programa se produce una colisión entre los nombres `<i>` y `<j>` de los patrones que ocurren en el patrón `<expr>` y en el patrón principal:

```
pl@nereida:~/Lperltesting$ cat -n clashofnamedofssets.pl
 1  #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1
 2  use v5.10;
 3
 4  my $input;
 5
 6  local $" = ", ";
 7
 8  my $parser = qr{
 9      ^ (?<i> (?&expr)) (?<j> (?&expr)) \z
10      (?{
11          say "main \${+} hash:";
12          say " (\$_ => \${+{\$_}}) " for sort keys %+;
13      })
14
15      (? (DEFINE)
16          (?<expr>
17              (?<i> . )
18              (?<j> . )
19              (?{
20                  say "expr \${+} hash:";
21                  say " (\$_ => \${+{\$_}}) " for sort keys %+;
22              })
23          )
24      )
25  }x;
26
27  $input = <>;
28  chomp($input);
29  if ($input =~ $parser) {
30      say "matches: ($&)";
31  }
```

La colisión hace que la salida sea esta:

```
pl@nereida:~/Lperltesting$ ./clashofnamedoffsets.pl
abab
expr $+ hash:
(i => a)
(j => b)
expr $+ hash:
(i => ab)
(j => b)
main $+ hash:
(i => ab)
(j => ab)
matches: (abab)
```

Si se evitan las colisiones, se evita la pérdida de información:

```
pl@nereida:~/Lperltesting$ cat -n namedoffsets.pl
 1  #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1
 2  use v5.10;
 3
 4  my $input;
 5
 6  local $" = ", ";
 7
 8  my $parser = qr{
 9      ^ (?<i> (?&expr)) (?<j> (?&expr)) \z
10      (?{
11          say "main \">$+ hash:";
12          say " ($_ => ${$_}) " for sort keys %+;
13      })
14
15      (? (DEFINE)
16          (?<expr>
17              (?<i_e> . )
18              (?<j_e> . )
19              (?{
20                  say "expr \">$+ hash:";
21                  say " ($_ => ${$_}) " for sort keys %+;
22              })
23          )
24      )
25  }x;
26
27  $input = <>;
28  chomp($input);
29  if ($input =~ $parser) {
30      say "matches: ($&)";
31  }
```

que al ejecutarse produce:

```
pl@nereida:~/Lperltesting$ ./namedoffsets.pl
abab
expr $+ hash:
```



```

(i_e => a)
(j_e => b)
expr $+ hash:
(i => ab)
(i_e => a)
(j_e => b)
main $+ hash:
(i => ab)
(j => ab)
matches: (abab)

```

3.2.9. Expresiones Regulares en tiempo de matching

Los paréntesis especiales:

```
(??{ Código Perl })
```

hacen que el Código Perl sea evaluado durante el tiempo de matching. El resultado de la evaluación se trata como una expresión regular. El match continuará intentando casar con la expresión regular retornada.

Paréntesis con memoria dentro de una *pattern code expression*

Los paréntesis en la expresión regular retornada no cuentan en el patrón exterior. Véase el siguiente ejemplo:

```

pl@nereida:~/Lperltesting$ cat -n postponedregexp.pl
1  #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w
2  use v5.10;
3  use strict;
4
5  my $r = qr{(?x)                # ignore spaces
6      ([ab])                    # save 'a' or 'b' in \$1
7      (??{ "($^N)"x3 })        # 3 more of the same as in \$1
8      };
9  say "<$&> lastpar = $#-" if 'bbbb' =~ $r;
10 say "<$&> lastpar = $#-" if 'aaaa' =~ $r;
11 say "<abab> didn't match" unless 'abab' =~ $r;
12 say "<aaab> didn't match" unless 'aaab' =~ $r;

```

Como se ve, hemos accedido desde el código interior al último paréntesis usando `$^N`. Sigue una ejecución:

```

pl@nereida:~/Lperltesting$ ./postponedregexp.pl
<bbbb> lastpar = 1
<aaaa> lastpar = 1
<abab> didn't match
<aaab> didn't match

```

Ejemplo: Secuencias de dígitos de longitud especificada por el primer dígito

Consideremos el problema de escribir una expresión regular que reconoce secuencias no vacías de dígitos tales que la longitud de la secuencia restante viene determinada por el primer dígito. Esta es una solución:

```

pl@nereida:~/Lperltesting$ cat -n intints.pl
1  #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w

```

```

2 use v5.10;
3 use strict;
4
5 my $r = qr{(?x)           # ignore spaces
6         (\d)             # a digit
7         ( ??{
8             "\\d{$^N}"    # as many as the former
9         })                # digit says
10        )
11    };
12 say "<$&> <$1> <$2>" if '3428' =~ $r;
13 say "<$&> <$1> <$2>" if '228' =~ $r;
14 say "<$&> <$1> <$2>" if '14' =~ $r;
15 say "24 does not match" unless '24' =~ $r;
16 say "4324 does not match" unless '4324' =~ $r;

```

Cuando se ejecuta se obtiene:

```

pl@nereida:~/Lperltesting$ ./intints.pl
<3428> <3> <428>
<228> <2> <28>
<14> <1> <4>
24 does not match
4324 does not match

```

Ejemplo: Secuencias de dígitos no repetidos

Otro ejemplo: queremos escribir una expresión regular que reconozca secuencias de n dígitos en las que no todos los dígitos se repiten. Donde quizá n es capturado de un paréntesis anterior en la expresión regular. Para simplificar la ilustración de la técnica supongamos que $n = 7$:

```

pl@nereida:~$ perl5.10.1 -wdE 0
main::(-e:1): 0
DB<1> x join '', map { "(?!".$_."{7})" } 0..9
0 '(?!0{7})(?!1{7})(?!2{7})(?!3{7})(?!4{7})(?!5{7})(?!6{7})(?!7{7})(?!8{7})(?!9{7})'
DB<2> x '7777777' =~ /(??{join '', map { "(?!".$_."{7})" } 0..9})(\d{7})/
empty array
DB<3> x '7777778' =~ /(??{join '', map { "(?!".$_."{7})" } 0..9})(\d{7})/
0 7777778
DB<4> x '4444444' =~ /(??{join '', map { "(?!".$_."{7})" } 0..9})(\d{7})/
empty array
DB<5> x '4422444' =~ /(??{join '', map { "(?!".$_."{7})" } 0..9})(\d{7})/
0 4422444

```

Palíndromos con independencia del acento

Se trata en este ejercicio de generalizar la expresión regular introducida en la sección 3.2.5 para reconocer los palabra-palíndromos.

Se trata de encontrar una regex que acepte que la lectura derecha e inversa de una frase en Español pueda diferir en la acentuación (como es el caso del clásico palíndromo *dábale arroz a la zorra el abad*). Una solución trivial es preprocesar la cadena eliminando los acentos. Supondremos sin embargo que se quiere trabajar sobre la cadena original. He aquí una solución:

```

1 pl@nereida:~/Lperltesting$ cat actionspanishpalin.pl
2 #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w -CIOEioA
3 use v5.10;

```

```

4 use strict;
5 use utf8;
6 use re 'eval';
7 use Switch;
8
9 sub f {
10     my $char = shift;
11
12     switch($char) {
13         case [ qw{a á} ] { return '[aá]' }
14         case [ qw{e é} ] { return '[eé]' }
15         case [ qw{i í} ] { return '[ií]' }
16         case [ qw{o ó} ] { return '[oó]' }
17         case [ qw{u ú} ] { return '[uú]' }
18         else { return $char };
19     }
20 }
21
22 my $regexp = qr/^(\\W* (?
23                 (\\w) (?-2)(??{ f($^N) })
24                 | \\w?
25                 ) \\W*
26             )
27 $
28 /ix;
29
30 my $input = <>; # Try: 'dábale arroz a la zorra el abad';
31 chomp($input);
32 if ($input =~ $regexp) {
33     say "$input is a palindrome";
34 }
35 else {
36     say "$input does not match";
37 }

```

Sigue un ejemplo de ejecución:

```

pl@nereida:~/Lperltesting$ ./actionspanishpalin.pl
dábale arroz a la zorra el abad
dábale arroz a la zorra el abad is a palindrome
pl@nereida:~/Lperltesting$ ./actionspanishpalin.pl
éoiúaáuióé
éoiúaáuióé is a palindrome
pl@nereida:~/Lperltesting$ ./actionspanishpalin.pl
dáed
dáed does not match

```

Postponiendo para conseguir recursividad

Véase el nodo Complex regex for maths formulas para la formulación del problema:

Hiya monks,

Im having trouble getting my head around a regular expression to match sequences. I need to catch all exceptions where a mathematical expression is illegal...

There must be either a letter or a digit either side of an operator parenthesis must open and close next to letters or digits, not next to operators, and do not have to exist variables must not be more than one letter Nothing other than a-z,A-Z,0-9,+,-,,/,,(,) can be used*

*Can anyone offer a hand on how best to tackle this problem?
many thanks*

La respuesta dada por ikegami usa (?{ ... }) para conseguir una conducta recursiva en versiones de perl anteriores a la 5.10:

```
pl@nereida:~/Lperltesting$ cat -n complexformula.pl
 1  #!/usr/bin/perl
 2  use strict;
 3  use warnings;
 4
 5  sub is_valid_expr {
 6      use re 'eval'; # to allow Eval-group at runtime
 7
 8      local our ($skip, $term, $expr);
 9      $skip = qr! \s* !x;
10      $term = qr! $skip [a-zA-Z]+          # A term is an identifier
11              | $skip [1-9][0-9]*          # or a number
12              | $skip \( (?{ $expr }) $skip # or an expression
13              \)
14              !x;
15      $expr = qr! $term                    # A expr is a term
16              (?: $skip [-+*/] $term )*    # or a term + a term ...
17              !x;
18
19      return $_[0] =~ / ^ $expr $skip \z /x;
20  }
21
22  print(is_valid_expr($_) ? "$_ is valid\n" : "$_ is not valid\n") foreach (
23      '(a + 3)',
24      '(3 * 4)+(b + x)',
25      '(5 - a)*z',
26      '3 + 2',
27
28      '!3 + 2',
29      '3 + 2!',
30
31      '3 a',
32      '3 3',
33      '3 * * 3',
34
35      '2 - 3 * 4',
36      '2 - 3 + 4',
37  );
```

Sigue el resultado de la ejecución:

```
pl@nereida:~/Lperltesting$ perl complexformula.pl
(a + 3) is valid
(3 * 4)+(b + x) is valid
```

```

(5 - a)*z is valid
3 + 2 is valid
!3 + 2 is not valid
3 + 2! is not valid
3 a is not valid
3 3 is not valid
3 * * 3 is not valid
2 - 3 * 4 is valid
2 - 3 + 4 is valid

```

Caveats

Estos son algunos puntos a tener en cuenta cuando se usan patrones postpuestos. Véase la entrada `(??{ code })` en la sección 'Extended-Patterns' en `perlre`:

WARNING: This extended regular expression feature is considered experimental, and may be changed without notice. Code executed that has side effects may not perform identically from version to version due to the effect of future optimisations in the regex engine.

This is a postponed regular subexpression. The code is evaluated at run time, at the moment this subexpression may match. The result of evaluation is considered as a regular expression and matched as if it were inserted instead of this construct.

The code is not interpolated.

As before, the rules to determine where the code ends are currently somewhat convoluted.

Because perl's regex engine is not currently re-entrant, delayed code may not invoke the regex engine either directly with `m//` or `s///`), or indirectly with functions such as `split`.

Recurring deeper than 50 times without consuming any input string will result in a fatal error. The maximum depth is compiled into perl, so changing it requires a custom build.

3.2.10. Expresiones Condicionales

Citando a `perlre`:

A conditional expression is a form of if-then-else statement that allows one to choose which patterns are to be matched, based on some condition.

There are two types of conditional expression: `(?(condition)yes-regexp)` and `(?(condition)yes-regexp)(?(condition)no-regexp)`. `(?(condition)yes-regexp)` is like an `if () { }` statement in Perl. If the condition is true, the `yes-regexp` will be matched. If the condition is false, the `yes-regexp` will be skipped and Perl will move onto the next `regexp` element.

The second form is like an `if () { } else { }` statement in Perl. If the condition is true, the `yes-regexp` will be matched, otherwise the `no-regexp` will be matched.

The condition can have several forms.

- *The first form is simply an integer in parentheses (integer). It is true if the corresponding backreference `\integer` matched earlier in the `regexp`. The same thing can be done with a name associated with a capture buffer, written as `(<name>)` or `('name')`.*
- *The second form is a bare zero width assertion `(?...)`, either a lookahead, a look-behind, or a code assertion.*
- *The third set of forms provides tests that return true if the expression is executed within a recursion `(R)` or is being called from some capturing group, referenced either by number `(R1)`, or by name `(R&name)`.*

Condiciones: número de paréntesis

Una expresión condicional puede adoptar diversas formas. La mas simple es un entero en paréntesis. Es cierta si la correspondiente referencia `\integer` casó (también se puede usar un nombre si se trata de un paréntesis con nombre).

En la expresión regular `/^(.)(..)?(?(2)a|b)/` si el segundo paréntesis casa, la cadena debe ir seguida de una `a`, si no casa deberá ir seguida de una `b`:

```
DB<1> x 'hola' =~ /^(.)(..)?(?(2)a|b)/
0 'h'
1 'ol'
DB<2> x 'ha' =~ /^(.)(..)?(?(2)a|b)/
empty array
DB<3> x 'hb' =~ /^(.)(..)?(?(2)a|b)/
0 'h'
1 undef
```

Ejemplo: cadenas de la forma *una-otra-otra-una*

La siguiente búsqueda casa con patrones de la forma `xx` o `xyyx`:

```
pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> x 'aa' =~ m{^(\\w+)(\\w+)?(?(2)\\2\\1|\\1)$}
0 'a'
1 undef
DB<2> x 'abba' =~ m{^(\\w+)(\\w+)?(?(2)\\2\\1|\\1)$}
0 'a'
1 'b'
DB<3> x 'abbc' =~ m{^(\\w+)(\\w+)?(?(2)\\2\\1|\\1)$}
empty array
DB<4> x 'juanpedropedrojuan' =~ m{^(\\w+)(\\w+)?(?(2)\\2\\1|\\1)$}
0 'juan'
1 'pedro'
```

Condiciones: Código

Una expresión condicional también puede ser un código:

```
DB<1> $a = 0; print "$&" if 'hola' =~ m{((?{$a})hola|adios)} # No hay matching

DB<2> $a = 1; print "$&" if 'hola' =~ m{((?{$a})hola|adios)}
hola
```

Ejemplo: Cadenas con posible paréntesis inicial (no anidados)

La siguiente expresión regular utiliza un condicional para forzar a que si una cadena comienza por un paréntesis abrir termina con un paréntesis cerrar. Si la cadena no comienza por paréntesis abrir no debe existir un paréntesis final de cierre:

```
pl@nereida:~/Lperltesting$ cat -n conditionalregexp.pl
1 #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1 -w
2  use v5.10;
3  use strict;
4
5  my $r = qr{(?x)                # ignore spaces
6      ^
7      ( \ ( )?                  # may be it comes an open par
```

```

8          [^()]+          # no parenthesis
9          (? (1)          # did we start with par?
10         \)              # if yes then close par
11        )
12       $
13      };
14     say "<$&>" if '(abcd)' =~ $r;
15     say "<$&>" if 'abc' =~ $r;
16     say "<(abc> does not match" unless '(abc' =~ $r;
17     say "<abc> does not match" unless 'abc)' =~ $r;

```

Al ejecutar este programa se obtiene:

```

pl@nereida:~/Lperltesting$ ./conditionalregexp.pl
<(abcd)>
<abc>
<(abc> does not match
<abc)> does not match

```

Expresiones Condicionales con (R)

El siguiente ejemplo muestra el uso de la condición (R), la cual comprueba si la expresión ha sido evaluada dentro de una recursión:

```

pl@nereida:~/Lperltesting$ perl5.10.1 -wdE 0
main::(-e:1): 0
DB<1> x 'bbaaaabb' =~ /(b(? (R) a+| (? 0)) b)/
0 'bbaaaabb'
DB<2> x 'bb' =~ /(b(? (R) a+| (? 0)) b)/
empty array
DB<3> x 'bab' =~ /(b(? (R) a+| (? 0)) b)/
empty array
DB<4> x 'bbabb' =~ /(b(? (R) a+| (? 0)) b)/
0 'bbabb'

```

La sub-expresión regular `(? (R) a+| (? 0))` dice: si esta siendo evaluada recursivamente admite **a+** si no, evalúa la regexp completa recursivamente.

Ejemplo: Palíndromos con Equivalencia de Acentos Españoles

Se trata en este ejercicio de generalizar la expresión regular introducida en la sección 3.2.5 para reconocer los palabra-palíndromos⁷. Se trata de encontrar una regexp que acepte que la lectura derecha e inversa de una frase en Español pueda diferir en la acentuación (como es el caso del clásico palíndromo *dábale arroz a la zorra el abad*). Una solución trivial es preprocesar la cadena eliminando los acentos. Supondremos sin embargo que se quiere trabajar sobre la cadena original. He aquí una solución parcial (por consideraciones de legibilidad sólo se consideran las vocales **a** y **o**:

```

1 pl@nereida:~/Lperltesting$ cat spanishpalin.pl
2 #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w -CIOEioA
3 use v5.10;
4 use strict;
5 use utf8;
6
7 my $regexp = qr/^(?<pal>\W* (? :
8                                     (?<L>( ?<a>[áa])|(?<e>[ée])|\w) # letter

```

⁷ No sé si existe el término. Significa que la lectura directa y la inversa pueden diferir en los signos de puntuación

```

9             (?&pal)                # nested palindrome
10            (?(<a>)[áa]             # if is an "a" group
11              | (?:((?<e>)[ée]      # if is an "e" group
12                |\g{L})             # exact match
13              )                     # end if [ée]
14            )                       # end group
15          )                         # end if [áa]
16        | \w?                       # non rec. case
17      ) \W*                         # punctuation symbols
18    )
19  $
20  /ix;
21
22  my $input = <>; # Try: 'dábale arroz a la zorra el abad';
23  chomp($input);
24  if ($input =~ $regex) {
25    say "$input is a palindrome";
26  }
27  else {
28    say "$input does not match";
29  }

```

Ejecución:

```

pl@nereida:~/Lperltesting$ ./spanishpalin.pl
dábale arroz a la zorra el abad
dábale arroz a la zorra el abad is a palindrome
pl@nereida:~/Lperltesting$ ./spanishpalin.pl
óuuo
óuuo does not match
pl@nereida:~/Lperltesting$ ./spanishpalin.pl
éaáé
éaáé is a palindrome

```

Hemos usado la opción `-CIOEioA` para asegurarnos que los ficheros de entrada/salida y error y la línea de comandos están en modo UTF-8. (Véase la sección ??)

Esto es lo que dice la documentación de `perlrun` al respecto:

The `-C` flag controls some of the Perl Unicode features.

As of 5.8.1, the `-C` can be followed either by a number or a list of option letters. The letters, their numeric values, and effects are as follows; listing the letters is equal to summing the numbers.

```

1  I 1 STDIN is assumed to be in UTF-8
2  O 2 STDOUT will be in UTF-8
3  E 4 STDERR will be in UTF-8
4  S 7 I + O + E
5  i 8 UTF-8 is the default PerlIO layer for input streams
6  o 16 UTF-8 is the default PerlIO layer for output streams
7  D 24 i + o
8  A 32 the @ARGV elements are expected to be strings encoded
9  in UTF-8
10 L 64 normally the "IOEioA" are unconditional,
11 the L makes them conditional on the locale environment

```



```

12 variables (the LC_ALL, LC_TYPE, and LANG, in the order
13 of decreasing precedence) -- if the variables indicate
14 UTF-8, then the selected "IOEioA" are in effect
15 a 256 Set ${^UTF8CACHE} to -1, to run the UTF-8 caching code in
16 debugging mode.

```

For example, -COE and -C6 will both turn on UTF-8-ness on both STDOUT and STDERR. Repeating letters is just redundant, not cumulative nor toggling.

The io options mean that any subsequent open() (or similar I/O operations) will have the :utf8 PerlIO layer implicitly applied to them, in other words, UTF-8 is expected from any input stream, and UTF-8 is produced to any output stream. This is just the default, with explicit layers in open() and with binmode() one can manipulate streams as usual.

-C on its own (not followed by any number or option list), or the empty string "" for the PERL_UNICODE environment variable, has the same effect as -CSDL . In other words, the standard I/O handles and the defaultopen() layer are UTF-8-fied but only if the locale environment variables indicate a UTF-8 locale. This behaviour follows the implicit (and problematic) UTF-8 behaviour of Perl 5.8.0.

You can use -CO (or 0 for PERL_UNICODE) to explicitly disable all the above Unicode features.

El pragma `use utf8` hace que se utilice una semántica de caracteres (por ejemplo, la regexp `/./` casará con un carácter unicode), el pragma `use bytes` cambia de semántica de caracteres a semántica de bytes (la regexp `.` casará con un byte).

3.2.11. Verbos que controlan el retroceso

El verbo de control (*FAIL)

Tomado de la sección 'Backtracking-control-verbs' en `perlretut`:

*The control verb (*FAIL) may be abbreviated as (*F). If this is inserted in a regexp it will cause to fail, just like at some mismatch between the pattern and the string. Processing of the regexp continues like after any "normal" failure, so that the next position in the string or another alternative will be tried. As failing to match doesn't preserve capture buffers or produce results, it may be necessary to use this in combination with embedded code.*

```

pl@nereida:~/Lperltesting$ cat -n vowelcount.pl
 1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w
 2  use strict;
 3
 4  my $input = shift() || <STDIN>;
 5  my %count = ();
 6  $input =~ /([aeiou])(?{ $count{$1}++; })(*FAIL)/i;
 7  printf("%s' => %3d\n", $_, $count{$_}) for (sort keys %count);

```

Al ejecutarse con entrada `supercalifragilistico` produce la salida:

```

pl@nereida:~/Lperltesting$ ./vowelcount.pl
supercalifragilistico
'a' =>    2
'e' =>    1
'i' =>    4
'o' =>    1
'u' =>    1

```

Ejercicio 3.2.5. ¿Que queda en `$1` después de ejecutado el matching `$input =~ /([aeiou])(?{ $count{$1}++; }`

Véase también:

- El nodo en PerlMonks *The Oldest Plays the Piano*
- Véase el ejercicio *Las tres hijas* en la sección 3.4.4

El verbo de control (*ACCEPT)

Tomado de perlretut:

*This pattern matches nothing and causes the end of successful matching at the point at which the (*ACCEPT) pattern was encountered, regardless of whether there is actually more to match in the string. When inside of a nested pattern, such as recursion, or in a subpattern dynamically generated via (??{}), only the innermost pattern is ended immediately.*

*If the (*ACCEPT) is inside of capturing buffers then the buffers are marked as ended at the point at which the (*ACCEPT) was encountered. For instance:*

```
DB<1> x 'AB' =~ /(A (A|B(*ACCEPT)|C) D)(E)/x
0  'AB'
1  'B'
2  undef
DB<2> x 'ACDE' =~ /(A (A|B(*ACCEPT)|C) D)(E)/x
0  'ACD'
1  'C'
2  'E'
```

El verbo SKIP

*This zero-width pattern prunes the backtracking tree at the current point when backtracked into on failure. Consider the pattern A (*SKIP) B, where A and B are complex patterns. Until the (*SKIP) verb is reached, A may backtrack as necessary to match. Once it is reached, matching continues in B, which may also backtrack as necessary; however, should B not match, then no further backtracking will take place, and the pattern will fail outright at the current starting position.*

*It also signifies that whatever text that was matched leading up to the (*SKIP) pattern being executed cannot be part of any match of this pattern. This effectively means that the regex engine skips forward to this position on failure and tries to match again, (assuming that there is sufficient room to match).*

*The name of the (*SKIP:NAME) pattern has special significance. If a (*MARK:NAME) was encountered while matching, then it is that position which is used as the "skip point". If no (*MARK) of that name was encountered, then the (*SKIP) operator has no effect. When used without a name the "skip point" is where the match point was when executing the (*SKIP) pattern.*

Ejemplo:

```
pl@nereida:~/Lperltesting$ cat -n SKIP.pl
1  #!/soft/perl5lib/bin/perl5.10.1 -w
2  use strict;
3  use v5.10;
4
5  say "NO SKIP: /a+b?(*FAIL)/";
6  our $count = 0;
7  'aaab' =~ /a+b?(?{print "$&\n"; $count++})(*FAIL)/;
8  say "Count=$count\n";
9
```

```

10 say "WITH SKIP: a+b?(*SKIP)(*FAIL)"/";
11 $count = 0;
12 'aaab' =~ /a+b?(*SKIP)(?{print "$&\n"; $count++})(*FAIL)/;
13 say "WITH SKIP: Count=$count\n";
14
15 say "WITH SKIP /a+(*SKIP)b?(*FAIL)/:";
16 $count = 0;
17 'aaab' =~ /a+(*SKIP)b?(?{print "$&\n"; $count++})(*FAIL)/;
18 say "Count=$count\n";
19
20 say "WITH SKIP /(*SKIP)a+b?(*FAIL): ";
21 $count = 0;
22 'aaab' =~ /(*SKIP)a+b?(?{print "$&\n"; $count++})(*FAIL)/;
23 say "Count=$count\n";

```

Ejecución:

```
pl@nereida:~/Lperltesting$ perl5.10.1 SKIP.pl
```

```
NO SKIP: /a+b?(*FAIL)/
```

```
aaab
```

```
aaa
```

```
aa
```

```
a
```

```
aab
```

```
aa
```

```
a
```

```
ab
```

```
a
```

```
Count=9
```

```
WITH SKIP: a+b?(*SKIP)(*FAIL)/
```

```
aaab
```

```
WITH SKIP: Count=1
```

```
WITH SKIP /a+(*SKIP)b?(*FAIL)/:
```

```
aaab
```

```
aaa
```

```
Count=2
```

```
WITH SKIP /(*SKIP)a+b?(*FAIL):
```

```
aaab
```

```
aaa
```

```
aa
```

```
a
```

```
aab
```

```
aa
```

```
a
```

```
ab
```

```
a
```

```
Count=9
```

Marcas

Tomado de la sección 'Backtracking-control-verbs' en `perlretut`:

(*MARK:NAME) (*:NAME)

*This zero-width pattern can be used to mark the point reached in a string when a certain part of the pattern has been successfully matched. This mark may be given a name. A later (*SKIP) pattern will then skip forward to that point if backtracked into on failure. Any number of (*MARK) patterns are allowed, and the NAME portion is optional and may be duplicated.*

*In addition to interacting with the (*SKIP) pattern, (*MARK:NAME) can be used to label a pattern branch, so that after matching, the program can determine which branches of the pattern were involved in the match.*

*When a match is successful, the \$REGMARK variable will be set to the name of the most recently executed (*MARK:NAME) that was involved in the match.*

This can be used to determine which branch of a pattern was matched without using a separate capture buffer for each branch, which in turn can result in a performance improvement.

*When a match has failed, and unless another verb has been involved in failing the match and has provided its own name to use, the \$REGERROR variable will be set to the name of the most recently executed (*MARK:NAME).*

```
pl@nereida:~/Lperltesting$ cat -n mark.pl
1 use v5.10;
2 use strict;
3
4 our $REGMARK;
5
6 $_ = shift;
7 say $REGMARK if /(?:x(*MARK:mx)|y(*MARK:my)|z(*MARK:mz))/;
8 say $REGMARK if /(?:x(*:xx)|y(*:yy)|z(*:zz))/;
```

Cuando se ejecuta produce:

```
pl@nereida:~/Lperltesting$ perl5.10.1 mark.pl y
my
yy
pl@nereida:~/Lperltesting$ perl5.10.1 mark.pl z
mz
zz
```

Poniendo un espacio después de cada signo de puntuación

Se quiere poner un espacio en blanco después de la aparición de cada coma:

`s/,/, /g;`

pero se quiere que la sustitución no tenga lugar si la coma esta incrustada entre dos dígitos. Además se pide que si hay ya un espacio después de la coma, no se duplique. Sigue una solución que usa marcas:

```
pl@nereida:~/Lperltesting$ perl5.10.1 -wde 0
main::(-e:1): 0
DB<1> $a = 'ab,cd, ef,12,34,efg,56,78,df, ef,'
DB<2> x ($b = $a) =~ s/\d,\d(*:d)|,(?!\\s)/($REGMARK eq 'd')? $& : ', '/ge
0 8
DB<3> p "<$b>"
<ab, cd, ef, 12,34, efg, 56,78, df, ef, >
```

3.3. Expresiones Regulares en Otros Lenguajes

Vim

- Learn vi/vim in 50 lines and 15 minutes
- VIM Regular Expressions
- Editing features for advanced users
- Vim documentation: pattern
- Vim Regular Expressions Chart

Java

El siguiente ejemplo muestra un programa estilo **grep**: solicita una expresión regular para aplicarla luego a una serie de entradas leídas desde la entrada estandar.

```
casiano@nereida:~/projects/PA/regexp$ cat -n Application.java
1  /**
2   * javac Application.java
3   * java Application
4   */
5
6  import java.io.*;
7  import java.util.regex.Pattern;
8  import java.util.regex.Matcher;
9
10 public class Application {
11
12     public static void main(String[] args){
13         String regexp = "";
14         BufferedReader br = new BufferedReader(new InputStreamReader(System.in));
15         try {
16             System.out.print("Enter your regex: ");
17             regexp = br.readLine();
18         } catch (IOException e) { System.exit(1); };
19         while (true) {
20
21             String input = "";
22             try {
23                 System.out.print("Enter input string to search: ");
24                 input = br.readLine();
25             } catch (IOException e) { System.exit(1); };
26
27             Pattern pattern = Pattern.compile(regexp);
28             Matcher matcher = pattern.matcher(input);
29
30             boolean found = false;
31             while (matcher.find()) {
32                 System.out.println("I found the text "
33                                     + matcher.group()
34                                     + " starting at index "
35                                     + matcher.start()
36                                     + " and ending at index "
```

```

37                                     +matcher.end()
38                                 );
39                                 found = true;
40                             }
41                             if(!found){
42                                 System.out.println("No match found.");
43                             }
44                         }
45                 }
46     }

```

Ejecución:

```

casiano@nereida:~/Ljavatesting$ java Application
Enter your regex: (\d+).(\d+)
Enter input string to search: a4b5d6c7efg
I found the text 4b5 starting at index 1 and ending at index 4
I found the text 6c7 starting at index 5 and ending at index 8
Enter input string to search: abc
No match found.
Enter input string to search:

```

Véase también Java Regular Expressions

bash

Esta es una versión en **bash** del conversor de temperaturas visto en las secciones anteriores:

```

pl@nereida:~/src/bash$ cat -n f2c
 1  #!/bin/bash
 2  echo "Enter a temperature (i.e. 32F, 100C):";
 3  read input;
 4
 5  if [ -z "$(echo $input | grep -i '^[+]\?[0-9]\+\(\.[0-9]*\)\?\ *[CF]$' )" ]
 6  then
 7      echo "Expecting a temperature, so don't understand \"$input\"." 1>&2;
 8  else
 9      input=$(echo $input | tr -d ' ');
10      InputNum=${input:0:${#input}-1};
11      Type=${input: -1}
12
13      if [ $Type = "c" -o $Type = "C" ]
14      then
15          celsius=$InputNum;
16          fahrenheit=$(echo "scale=2; ($celsius * 9/5)+32" | bc -l);
17      else
18          fahrenheit=$InputNum;
19          celsius=$(echo "scale=2; ($fahrenheit -32)*5/9" | bc -l);
20      fi
21
22      echo "$celsius C = $fahrenheit F";
23  fi

```

C

```

pl@nereida:~/src/regexpr$ cat -n pcregrep.c
 1  #include <stdio.h>
 2  #include <stdlib.h>
 3  #include <string.h>
 4  #include <assert.h>
 5  #include <pcre.h>
 6
 7  char enter_reverse_mode[] = "\33[7m";
 8  char exit_reverse_mode[] = "\33[0m";
 9
10  int main(int argc, char **argv)
11  {
12      const char *pattern;
13      const char *errstr;
14      int erroffset;
15      pcre *expr;
16      char line[512];
17      assert(argc == 2); /* XXX fixme */
18      pattern = argv[1];
19      if (!(expr = pcre_compile(pattern, 0, &errstr, &erroffset, 0))) {
20          fprintf(stderr, "%s: %s\n", pattern, errstr);
21          return EXIT_FAILURE;
22      }
23      while (fgets(line, sizeof line, stdin)) {
24          size_t len = strcspn(line, "\n");
25          int matches[2];
26          int offset = 0;
27          int flags = 0;
28          line[len] = '\0';
29          while (0 < pcre_exec(expr, 0, line, len, offset, flags, matches, 2)) {
30              printf("%. *s%s%. *s%s",
31                  matches[0] - offset, line + offset,
32                  enter_reverse_mode,
33                  matches[1] - matches[0], line + matches[0],
34                  exit_reverse_mode);
35              offset = matches[1];
36              flags |= PCRE_NOTBOL;
37          }
38          printf("%s\n", line + offset);
39      }
40      return EXIT_SUCCESS;
41  }

```

Compilación:

```
pl@nereida:~/src/regexpr$ gcc -lpcre pcregrep.c -o pcregrep
```

Cuando se ejecuta espera un patrón en la línea de comandos y pasa a leer desde la entrada estandar. Las cadenas que casan se muestran resaltadas:

```
pl@nereida:~/src/regexpr$ ./pcregrep '\d+'
```

```
435 otro 23
```

```
435 otro 23
```

```
hola
```

```
hola
```

Python

```
pl@nereida:~/src/python$ cat -n c2f.py
1  #!/usr/local/bin/python
2  import re
3
4  temp = raw_input( ' Introduzca una temperatura (i.e. 32F, 100C): ' )
5  pattern = re.compile( "^[+-]?[0-9]+(\\.[0-9]*)?\\s*([CF])$", re.IGNORECASE )
6  mo = pattern.match( temp )
7
8  if mo:
9      inputNum = float(mo.group( 1 ))
10     type = mo.group( 3 )
11     celsius = 0.0
12     fahrenheit = 0.0
13     if ( type == "C" or type == "c" ) :
14         celsius = inputNum
15         fahrenheit = ( celsius * 9/5 ) + 32
16     else :
17         fahrenheit = inputNum
18         celsius = ( fahrenheit - 32 ) * 5/9
19     print " ", '%.2f'%(celsius), " C = ", '%.2f'%(fahrenheit), " F\n"
20 else :
21     print " Se esperaba una temperatura, no se entiende", temp, "\n"
```

Ruby

```
pl@nereida:~/src/ruby$ cat -n f2c_b
1  #!/usr/bin/ruby
2
3  # Primero leemos una temperatura
4  class Temperature_calculator
5      def initialize temp
6          comp = Regexp.new('^[+-]?\\d+(\\.\\d*)?\\s*([CFcf])$')
7          if temp =~ comp
8              begin
9                  cifra = Float($1)
10                 @C,@F = ( $3 == "F" or $3 == "f"? [(cifra -32) * 5/9, cifra] : [cifra , cifra * 9/5 +
11             end
12             else
13                 raise("Entrada incorrecta")
14             end
15         end
16
17         def show
18             puts "Temperatura en Celsius: #{@C}, temperatura en Fahrenheit: #{@F}"
19         end
20     end
21
22     temperatura = Temperature_calculator.new(readline.chop)
23     temperatura.show
```

Javascript


```

<SCRIPT LANGUAGE="JavaScript"><!--
function demoMatchClick() {
    var re = new RegExp(document.demoMatch.regex.value);
    if (document.demoMatch.subject.value.match(re)) {
        alert("Successful match");
    } else {
        alert("No match");
    }
}

function demoShowMatchClick() {
    var re = new RegExp(document.demoMatch.regex.value);
    var m = re.exec(document.demoMatch.subject.value);
    if (m == null) {
        alert("No match");
    } else {
        var s = "Match at position " + m.index + ":\n";
        for (i = 0; i < m.length; i++) {
            s = s + m[i] + "\n";
        }
        alert(s);
    }
}

function demoReplaceClick() {
    var re = new RegExp(document.demoMatch.regex.value, "g");
    document.demoMatch.result.value =
        document.demoMatch.subject.value.replace(re,
            document.demoMatch.replacement.value);
}
// -->
</SCRIPT>

<FORM ID="demoMatch" NAME="demoMatch" METHOD=POST ACTION="javascript:void(0)">
<P>Regex: <INPUT TYPE=TEXT NAME="regex" VALUE="\bt[a-z]+\b" SIZE=50></P>
<P>Subject string: <INPUT TYPE=TEXT NAME="subject"
    VALUE="This is a test of the JavaScript RegExp object" SIZE=50></P>
<P><INPUT TYPE=SUBMIT VALUE="Test Match" ONCLICK="demoMatchClick()">
<INPUT TYPE=SUBMIT VALUE="Show Match" ONCLICK="demoShowMatchClick()"></P>

<P>Replacement text: <INPUT TYPE=TEXT NAME="replacement" VALUE="replaced" SIZE=50></P>
<P>Result: <INPUT TYPE=TEXT NAME="result"
    VALUE="click the button to see the result" SIZE=50></P>
<P><INPUT TYPE=SUBMIT VALUE="Replace" ONCLICK="demoReplaceClick()"></P>
</FORM>

```

3.4. Casos de Estudio

3.4.1. Secuencias de números de tamaño fijo

El siguiente problema y sus soluciones se describen en el libro de J.E.F. Friedl [?]. Supongamos que tenemos un texto conteniendo códigos que son números de tamaño fijo, digamos seis dígitos, todos pegados, sin separadores entre ellos, como sigue:

012345678901**123334**2345678901231**25934**890123345126

El problema es encontrar los códigos que comienzan por 12. En negrita se han resaltado las soluciones. Son soluciones sólo aquellas que, comienzan por 12 en una posición múltiplo de seis. Una solución es:

```
@nums = grep {m/^12/} m/\d{6}/g;
```

que genera una lista con los números y luego selecciona los que comienzan por 12. Otra solución es:

```
@nums = grep { defined } m/(12\d{4})|\d{6}/g;
```

que aprovecha que la expresión regular devolverá una lista vacía cuando el número no empieza por 12:

```
DB<1> $x = '012345678901123334234567890123125934890123345126'
```

```
DB<2> x ($x =~ m/(12\d{4})|\d{6}/g)
```

```
0 undef
1 undef
2 123334
3 undef
4 undef
5 125934
6 undef
7 undef
```

Obsérvese que se está utilizando también que el operador `|` no es *greedy*.

¿Se puede resolver el problema usando sólo una expresión regular? Obsérvese que esta solución “casi funciona”:

```
DB<3> x @nums = $x =~ m/(?:\d{6})*?(12\d{4})/g;
```

```
0 123334
1 125934
2 123345
```

recoge la secuencia mas corta de grupos de seis dígitos que no casan, seguida de una secuencia que casa. El problema que tiene esta solución es al final, cuando se han casado todas las soluciones, entonces la búsqueda exhaustiva hará que nos muestre soluciones que no comienzan en posiciones múltiplo de seis. Por eso encuentra 123345:

012345678901**123334**2345678901231**25934**8901**23345**126

Por eso, Friedl propone esta solución:

```
@nums = m/(?:\d{6})*?(12\d{4})(?:(!12)\d{6})*/g;
```

Se asume que existe al menos un éxito en la entrada inicial. Que es un extraordinario ejemplo de como el uso de paréntesis de agrupamiento simplifica y mejora la legibilidad de la solución. Es fantástico también el uso del operador de predicción negativo.

Solución usando el ancla \ G

El ancla \G ha sido concebida para su uso con la opción /g. Casa con el punto en la cadena en el que terminó el último emparejamiento. Cuando se trata del primer intento o no se está usando /g, usar \G es lo mismo que usar \A.

Mediante el uso de este ancla es posible formular la siguiente solución al problema planteado:

```
pl@nereida:~/Lperltesting$ perl -wde 0
main::(-e:1): 0
DB<1> $_ = '012345678901123334234567890123125934890123345126'
DB<2> x m/\G(?:\d{6})*?(12\d{4})/g
0 123334
1 125934
```

Sustitución

Si lo que se quiere es sustituir las secuencias deseadas es posible hacerlo con la siguiente expresión regular:

```
casiano@nereida:~/docs/curriculums/CV_MEC$ perl -wde 0
DB<1> x $x = '012345678901123334234567890123125934890123345126'
0 012345678901123334234567890123125934890123345126
DB<2> x ($y = $x) =~ s/(12\d{4})|\d{6}/$1? "-$1-":$& /ge
0 8
DB<3> p $y
012345678901-123334-234567890123-125934-890123345126
```

3.4.2. Palabras Repetidas

Su jefe le pide una herramienta que compruebe la aparición de duplicaciones consecutivas en un texto texto (como esta esta y la anterior anterior). La solución debe cumplir las siguientes especificaciones:

- Aceptar cualquier número de ficheros. Resaltar las apariciones de duplicaciones. Cada línea del informe debe estar precedida del nombre del fichero.
- Funcionar no sólo cuando la duplicación ocurre en la misma línea.
- Funcionar independientemente del *case* y de los blancos usados en medio de ambas palabras.
- Las palabras en cuestión pueden estar separadas por *tags* HTML.

```
1 #!/usr/bin/perl -w
2 use strict;
3 use Term::ANSIScreen qw/:constants/;
4
5 my $bold = BOLD();
6 my $clear = CLEAR();
7 my $line = 1;
8
9 # read paragraph
10 local $/ = ".\n";
11 while (my $par = <>) {
12     next unless $par =~ s{
13         \b                # start word ...
14         ([a-z]+)          # grab word in $1 and \1
15         (                  # save the tags and spaces in $2
16         (\s|<[^>]+>)+     # spaces or HTML tags
```

```

17         )
18         (\1\b)           # repeated word in $4
19     }!$bold$1$clear$2$bold$4$clear!igx;
20
21     $par =~ s/~/ "$ARGV("$line++")": "/meg;    # insert filename and line number
22
23     print $par;
24 }

```

3.4.3. Análisis de cadenas con datos separados por comas

Supongamos que tenemos cierto texto en `$text` proveniente de un fichero CSV (*Comma Separated Values*). Esto es el fichero contiene líneas con el formato:

```
"earth",1,"moon",9.374
```

Esta línea representa cinco campos. Es razonable querer guardar esta información en un *array*, digamos `@field`, de manera que `$field[0] == 'earth'`, `$field[1] == '1'`, etc. Esto no sólo implica descomponer la cadena en campos sino también quitar las comillas de los campos entrecomillados. La primera solución que se nos ocurre es hacer uso de la función `split`:

```
@fields = split(/,/, $text);
```

Pero esta solución deja las comillas dobles en los campos entrecomillados. Peor aún, los campos entrecomillados pueden contener comas, en cuyo caso la división proporcionada por `split` sería errónea.

```

1  #!/usr/bin/perl -w
2  use Text::ParseWords;
3
4  sub parse_csv {
5      my $text = shift;
6      my @fields = (); # initialize @fields to be empty
7
8      while ($text =~
9          m/"((["\\]|\\.)*")/? # quoted fields
10         |
11         ([^,]+),?           # $3 = non quoted fields
12         |
13         ,                   # allows empty fields
14         /gx
15     )
16     {
17         push(@fields, defined($1)? $1:$3); # add the just matched field
18     }
19     push(@fields, undef) if $text =~ m/,$/; #account for an empty last field
20     return @fields;
21 }
22
23 $test = '"earth",1,"a1, a2","moon",9.374';
24 print "string = \"$test\"";
25 print "Using parse_csv\n";
26 @fields = parse_csv($test);
27 foreach $i (@fields) {
28     print "$i\n";

```

```

29 }
30
31 print "Using Text::ParseWords\n:";
32 # @words = &quotewords($delim, $keep, @lines);
33 #The $keep argument is a boolean flag. If true, then the
34 #tokens are split on the specified delimiter, but all other
35 #characters (quotes, backslashes, etc.) are kept in the
36 #tokens. If $keep is false then the &*quotewords()
37 #functions remove all quotes and backslashes that are not
38 #themselves backslash-escaped or inside of single quotes
39 #(i.e., &quotewords() tries to interpret these characters
40 #just like the Bourne shell).
41
42 @fields = quotewords(',',',0,$test);
43 foreach $i (@fields) {
44     print "$i\n";
45 }

```

Las subrutinas en Perl reciben sus argumentos en el *array* `@_`. Si la lista de argumentos contiene listas, estas son “aplanadas” en una única lista. Si, como es el caso, la subrutina ha sido declarada antes de la llamada, los argumentos pueden escribirse sin paréntesis que les rodeen:

```
@fields = parse_csv $test;
```

Otro modo de llamar una subrutina es usando el prefijo `&`, pero sin proporcionar lista de argumentos.

```
@fields = &parse_csv;
```

En este caso se le pasa a la rutina el valor actual del *array* `@_`.

Los operadores `push` (usado en la línea 17) y `pop` trabajan sobre el final del *array*. De manera análoga los operadores `shift` y `unshift` lo hacen sobre el comienzo. El operador ternario `?` trabaja de manera análoga como lo hace en C.

El código del `push` podría sustituirse por este otro:

```
push(@fields, $+);
```

Puesto que la variable `$+` contiene la cadena que ha casado con el último paréntesis que haya casado en el último “matching”.

La segunda parte del código muestra que existe un módulo en Perl, el módulo `Text::Parsewords` que proporciona la rutina `quotewords` que hace la misma función que nuestra subrutina.

Sigue un ejemplo de ejecución:

```

> csv.pl
string = ' "earth",1,"a1, a2","moon",9.374'
Using parse_csv
:earth
1
a1, a2
moon
9.374
Using Text::ParseWords
:earth
1
a1, a2
moon
9.374

```

3.4.4. Las Expresiones Regulares como Exploradores de un Árbol de Soluciones

Números Primos

El siguiente programa evalúa si un número es primo o no:

```
pl@nereida:~/Lperltesting$ cat -n isprime.pl
1  #!/usr/bin/perl -w
2  use strict;
3
4  my $num = shift;
5  die "Usage: $0 integer\n" unless (defined($num) && $num =~ /\d+$/);
6
7  if (("1" x $num) =~ /^(11+)\1+$/) {
8      my $factor = length($1);
9      print "$num is $factor x ".$num/$factor."\n";
10 }
11 else {
12     print "$num is prime\n";
13 }
```

Siguen varias ejecuciones:

```
pl@nereida:~/Lperltesting$ ./isprime.pl 35.32
Usage: ./isprime.pl integer
pl@nereida:~/Lperltesting$ ./isprime.pl 47
47 is prime
pl@nereida:~/Lperltesting$ ./isprime.pl 137
137 is prime
pl@nereida:~/Lperltesting$ ./isprime.pl 147
147 is 49 x 3
pl@nereida:~/Lperltesting$ ./isprime.pl 137
137 is prime
pl@nereida:~/Lperltesting$ ./isprime.pl 49
49 is 7 x 7
pl@nereida:~/Lperltesting$ ./isprime.pl 47
47 is prime
```

Ecuaciones Diofánticas: Una solución

Según dice la entrada *Diophantine equation* en la wikipedia:

In mathematics, a Diophantine equation is an indeterminate polynomial equation that allows the variables to be integers only.

La siguiente sesión con el depurador muestra como se puede resolver una ecuación lineal diofántica con coeficientes positivos usando una expresión regular:

```
DB<1> # Resolvamos 3x + 2y + 5z = 40
DB<2> x ('a'x40) =~ /^(?:(...)+)((?:...)+)((?:.....)+)$/
0  'aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa'
1  'aa'
2  'aaaaa'
DB<3> x map { length } ('a'x40) =~ /^(?:(...)+)((?:...)+)((?:.....)+)$/
0  33
1  2
2  5
```

```
DB<4> @c = (3, 2, 5)
DB<5> x map { length($_) / $c[$i++] } ('a'x40) =~ /^(?:(...)+)((?:...)+)((?:.....)+)$/
0 11
1 1
2 1
DB<6> p 3*11+2*1+5*1
40
```

Ecuaciones Diofánticas: Todas las soluciones

Usando el verbo (*FAIL) es posible obtener todas las soluciones:

```
main::(-e:1): 0
DB<1> sub equ { my @c = @_; print "\t3*$c[0]+2*$c[1]+5*$c[2] = ", 3*$c[0]+2*$c[1]+5*$c[2], "\n"
DB<2> sub f { my @c = ((length($1)/3), (length($2)/2), (length($3)/5)); equ(@c); }
DB<3> x ('a'x40) =~ /^(?:(...)+)((?:...)+)((?:.....)+)$({ f() })(*FAIL)/x
    3*11+2*1+5*1 = 40
    3*9+2*4+5*1 = 40
    3*8+2*3+5*2 = 40
    3*7+2*7+5*1 = 40
    3*7+2*2+5*3 = 40
    3*6+2*6+5*2 = 40
    3*6+2*1+5*4 = 40
    3*5+2*10+5*1 = 40
    3*5+2*5+5*3 = 40
    3*4+2*9+5*2 = 40
    3*4+2*4+5*4 = 40
    3*3+2*13+5*1 = 40
    3*3+2*8+5*3 = 40
    3*3+2*3+5*5 = 40
    3*2+2*12+5*2 = 40
    3*2+2*7+5*4 = 40
    3*2+2*2+5*6 = 40
    3*1+2*16+5*1 = 40
    3*1+2*11+5*3 = 40
    3*1+2*6+5*5 = 40
    3*1+2*1+5*7 = 40
    empty array
DB<4>
```

Ecuaciones Diofánticas: Resolutor general

El siguiente programa recibe en línea de comandos los coeficientes y término independiente de una ecuación lineal diofántica con coeficientes positivos y muestra todas las soluciones. El algoritmo primero crea una cadena conteniendo el código Perl que contiene la expresión regular adecuada para pasar luego a evaluarlo:

```
pl@nereida:~/Lperltesting$ cat -n diophantinesolvergen.pl
1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w
2  use v5.10;
3  use strict;
4
5  # Writes a Perl solver for
6  # a1 x1 + a2 x2 + ... + an xn = b
7  # a_i and b integers > 0
8  #
```



```

6 6 2
6 1 4
5 10 1
5 5 3
4 9 2
4 4 4
3 13 1
3 8 3
3 3 5
2 12 2
2 7 4
2 2 6
1 16 1
1 11 3
1 6 5
1 1 7

```

Las Tres Hijas

En la páginas de Retos Matemáticos de
DIVULGAMAT
puede encontrarse el siguiente problema:

Ejercicio 3.4.1. *Dos matemáticos se vieron en la calle después de muchos años sin coincidir.*

- *¡Hola!, ¿qué tal?, ¿te casaste?, y... ¿cuántos hijos tienes?*
- *Pues tengo tres hijas.*
- *¿y qué años tienen?*
- *¡A ver si lo adivinas!: el producto de las edades de las tres es 36, y su suma es el número del portal que ves enfrente...*
- *¡Me falta un dato!*
- *¡Ah, sí!, ¡la mayor toca el piano!*

¿Qué edad tendrán las tres hijas?

¿Podemos ayudarnos de una expresión regular para resolver el problema? Al ejecutar el siguiente programa:

```

pl@nereida:~/Lperltesting$ cat -n playspiano.pl
1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1 -w
2  use v5.10;
3  use strict;
4  use List::Util qw{sum};
5
6  local our %u;
7  sub f {
8      my @a = @_;
9      @a = sort { $b <=> $a } (length($a[1]), length($a[0])/length($a[1]), 36/length($a[0]) );
10
11     local $" = ", ";
12     say "(@a)\t ".sum(@a) unless exists($u{"@a"});
13     $u{"@a"} = undef;
14 }

```

```

15
16 say "SOL\t\tNUMBER";
17 my @a = ('1'x36) =~
18         /\^((1+)\2+)(\1+)$
19         (?{ f($1, $2, $3)
20             })
21         (*FAIL)
22         /x;

```

obtenemos la salida:

```

pl@nereida:~/Lperltesting$ ./playspiano.pl
SOL          NUMBER
(9, 2, 2)      13
(6, 3, 2)      11
(4, 3, 3)      10
(18, 2, 1)     21
(12, 3, 1)     16
(9, 4, 1)      14
(6, 6, 1)      13

```

Explique el funcionamiento del programa. A la vista de la salida ¿Cuáles eran las edades de las hijas?

Mochila 0-1

Para una definición del problema vea la sección El Problema de la Mochila 0-1 en los apuntes de LHP

Ejercicio 3.4.2. *¿Sería capaz de resolver usando expresiones regulares el problema de la mochila 0-1? ¿Si lo logra merece el premio a la solución mas freak que se haya encontrado para dicho problema!*

Véase también

Véase también:

- Véase el nodo en PerlMonks *The Oldest Plays the Piano*
- Solving Algebraic Equations Using Regular Expressions

3.4.5. Número de substituciones realizadas

El operador de substitución devuelve el número de substituciones realizadas, que puede ser mayor que uno si se usa la opción /g. En cualquier otro caso retorna el valor falso.

```

1 #!/usr/bin/perl -w
2 undef($/);
3 $paragraph = <STDIN>;
4 $count = 0;
5 $count = ($paragraph =~ s/Mister\b/Mr./ig);
6 print "$paragraph";
7 print "\n$count\n";

```

El resultado de la ejecución es el siguiente:

```

> numsust.pl
Dear Mister Bean,
Is a pleasure for me and Mister Pluto

```

to invite you to the Opening Session
Official dinner that will be chaired by
Mister Goofy.

Yours sincerely

Mister Mickey Mouse

Dear Mr. Bean,

Is a pleasure for me and Mr. Pluto
to invite you to the Opening Session
Official dinner that will be chaired by
Mr. Goofy.

Yours sincerely

Mr. Mickey Mouse

4

3.4.6. Expandiendo y comprimiendo tabs

Este programa convierte los tabs en el número apropiado de blancos.

```
pl@nereida:~/Lperltesting$ cat -n expandtabs.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3
 4  my @string = <>;
 5
 6  for (@string) {
 7      while (s/\t+/' ' x (length($&)*8 - length($')%8)/e) {}
 8      print $_;
 9  }
```

Sigue un ejemplo de ejecución:

```
pl@nereida:~/Lperltesting$ cat -nt tabs.in
 1  012345670123456701234567012345670
 2  one^Itwo^I^Ithree
 3  four^I^I^I^Ifive
 4  ^I^Itwo

pl@nereida:~/Lperltesting$ ./expandtabs.pl tabs.in | cat -tn
 1  012345670123456701234567012345670
 2  one      two      three
 3  four                                five
 4                                two
```

Ejercicio 3.4.3. *¿Funciona igual si se cambia el bucle while por una opción /g?*

```
pl@nereida:~/Lperltesting$ cat -n ./expandtabs2.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3
 4  my @string = <>;
 5
 6  for (@string) {
 7      s/\t+/' ' x (length($&)*8 - length($')%8)/ge;
 8      print $_;
 9  }
```

¿Porqué?

3.4.7. Modificación de Múltiples Ficheros: one liner

Aunque no es la forma de uso habitual, Perl puede ser utilizado en “modo sed” para modificar el texto en múltiples ficheros:

```
perl -e 's/nereida\.deioc\.ull\.es/miranda.deioc.ull.es/gi' -p -i.bak *.html
```

Este programa sustituye la palabra original (g)lobalmente e i)gnorando el “case”) en todos los ficheros *.html y para cada uno de ellos crea una copia de seguridad *.html.bak.

Otro ejemplo: la sustitución que sigue ocurre en todos los ficheros info.txt en todos los subdirectorios de los subdirectorios que comiencen por alu:

```
perl -e 's/\\|hyperpage//gi' -p -i.bak alu*/*/info.txt
```

Las *opciones de línea* de comandos significan lo siguiente:

-e puede usarse para definir el script en la línea de comandos. Múltiples -e te permiten escribir un multi-script. Cuando se usa -e, perl no busca por un fichero de script entre la lista de argumentos.

-p La opción -p hace que perl incluya un bucle alrededor de tu “script” al estilo sed:

```
while (<>) {  
    ...                # your script goes here  
} continue {  
    print;  
}
```

-n Nótese que las líneas se imprimen automáticamente. Para suprimir la impresión usa la opción -n

-i[ext] La opción -i Expresa que los ficheros procesados serán modificados. Se renombra el fichero de entrada file.in a file.in.ext, abriendo el de salida con el mismo nombre del fichero de entrada file.in. Se selecciona dicho fichero como de salida por defecto para las sentencias print. Si se proporciona una extensión se hace una copia de seguridad. Si no, no se hace copia de seguridad.

En general las opciones pueden ponerse en la primera línea del “script”, donde se indica el intérprete. Así pues, decir

```
perl -p -i.bak -e "s/foo/bar/;"  
es equivalente a usar el “script”:
```

```
#!/usr/bin/perl -pi.bak  
s/foo/bar/;
```

3.5. tr y split

El operador de traducción permite la conversión de unos caracteres por otros. Tiene la sintaxis:

```
tr/SEARCHLIST/REPLACEMENTLIST/cds  
y/SEARCHLIST/REPLACEMENTLIST/cds
```

El operador permite el reemplazo carácter a carácter, por ejemplo:

```
$ perl -de 0  
DB<1> $a = 'fiboncacci'  
DB<2> $a =~ tr/aeiou/AEIOU/  
DB<3> print $a  
fIb0ncAccI
```

```
DB<4> $a =~ y/fbnc/FBNC/
DB<5> print $a
FIBONCACCI
```

El operador devuelve el número de caracteres reemplazados o suprimidos.

```
$cnt = $sky =~ tr/*/*/; # count the stars in $sky
```

Si se especifica el modificador /d, cualquier carácter en SEARCHLIST que no figure en REPLACEMENTLIST es eliminado.

```
DB<6> print $a
FIBONCACCI
DB<7> $a =~ y/OA//d
DB<8> print $a
FIBNCCCI
```

Si se especifica el modificador /s, las secuencias de caracteres consecutivos que serían traducidas al mismo carácter son comprimidas a una sola:

```
DB<1> $b = 'aaghhh!'
DB<2> $b =~ tr/ah//s
DB<3> p $b
agh!
```

Observa que si la cadena REPLACEMENTLIST es vacía, no se introduce ninguna modificación.

Si se especifica el modificador /c, se complementa SEARCHLIST; esto es, se buscan los caracteres que no están en SEARCHLIST.

```
tr/a-zA-Z/ /cs; # change non-alphas to single space
```

Cuando se dan múltiples traducciones para un mismo carácter, solo la primera es utilizada:

```
tr/AAA/XYZ/
```

traducirá A por X.

El siguiente *script* busca una expresión regular en el fichero de *passwords* e imprime los *login* de los usuarios que casan con dicha cadena. Para evitar posibles confusiones con las vocales acentuadas se usa el operador tr.

```
1 #!/usr/bin/perl -w
2 $search = shift(@ARGV) or die("you must provide a regexpr\n");
3 $search =~ y/ÁÉÍÓÚáéíóú/AEIOUaeiou/;
4 open(FILE, "/etc/passwd");
5 while ($line = <FILE>) {
6     $line =~ y/ÁÉÍÓÚáéíóú/AEIOUaeiou/;
7     if ($line =~ /$search/io) {
8         @fields = split(":", $line);
9         $login = $fields[0];
10        if ($line !~ /^#/ ) {
11            print "$login\n";
12        }
13        else {
14            print "#$login\n";
15        }
16    }
17 }
18
```

Ejecución (suponemos que el nombre del fichero anterior es `split.pl`):

```
> split.pl Rodriguez
##direccion
call
casiano
alu5
alu6
##doctorado
paco
falmeida
##ihiu07
```

Para familiarizarte con este operador, codifica y prueba el siguiente código:

```
1 #!/usr/bin/perl -w
2 $searchlist = shift @ARGV;
3 $replacelist = shift @ARGV;
4 $option = "";
5 $option = shift @ARGV if @ARGV;
6
7 while (<>) {
8     $num = eval "tr/$searchlist/$replacelist/$option";
9     die "$@" if $@;
10    print "$num: $_";
11 }
```

Perl construye la tabla de traducción en “tiempo de compilación”. Por ello ni `SEARCHLIST` ni `REPLACEMENTLIST` son susceptibles de ser interpolados. Esto significa que si queremos usar variables tenemos que recurrir a la función `eval`.

La expresión pasada como parámetro a `eval` en la línea 8 es analizada y ejecutada como si se tratara de un pequeño programa Perl. Cualquier asignación a variables permanece después del `eval`, así como cualquier definición de subrutina. El código dentro de `eval` se trata como si fuera un bloque, de manera que cualesquiera variables locales (declaradas con `my`) desaparecen al final del bloque.

La variable `$@` contiene el mensaje de error asociado con la última ejecución del comando `eval`. Si es nula es que el último comando se ejecuto correctamente. Aquí tienes un ejemplo de llamada:

```
> tr.pl 'a-z' 'A-Z' s
jose hernandez
13: JOSE HERNANDEZ
joosee hernnandez
16: JOSE HERNANDEZ
```

3.6. Pack y Unpack

El operador `pack` trabaja de forma parecida a `sprintf`. Su primer argumento es una cadena, seguida de una lista de valores a formatear y devuelve una cadena:

```
pack("CCC", 65, 66, 67, 68) # empaquetamos A B C D
```

el inverso es el operador `unpack`

```
unpack("CCC", "ABCD")
```

La cadena de formato es una lista de especificadores que indican el tipo del dato que se va a empaquetar/desempaquetar. Cada especificador puede opcionalmente seguirse de un contador de repetición que indica el número de elementos a formatear. Si se pone un asterisco (*) se indica que la especificación se aplica a todos los elementos restantes de la lista.

Formato	Descripción
A	Una cadena completada con blancos
a	Una cadena completada con ceros
B	Una cadena binaria en orden descendente
b	Una cadena binaria en orden ascendente
H	Una cadena hexadecimal, los nibble altos primero
h	Una cadena hexadecimal, los nibble bajos primero

Ejemplo de uso del formato A:

```
DB<1> $a = pack "A2A3", "Pea","r1"
DB<2> p $a
Perl
DB<3> @b = unpack "A2A3", "Perl"
DB<4> p "@b"
Pe r1
```

La variable @b tiene ahora dos cadenas. Una es Pe la otra es r1. Veamos un ejemplo con el formato B:

```
p ord('A')
65
DB<22> $x = pack "B8", "01000001"
DB<23> p $x
A
DB<24> @y = unpack "B8", "A"
DB<25> p "@y"
01000001
DB<26> $x = pack "b8", "10000010"
DB<27> p $x
```

3.7. Práctica: Un lenguaje para Componer Invitaciones

En el capítulo 6 (sección 6.4.2.2) del libro The LaTeX Web Companion se define un lenguaje para componer textos para enviar invitaciones.

Para escribir una invitación en ese lenguaje escribiríamos algo así:

```
pl@nereida:~/Lp10910/Practicas/161009/src$ cat -n invitation.xml
1  <?xml version="1.0"?>
2  <!DOCTYPE invitation SYSTEM "invitation.dtd">
3  <invitation>
4  <!-- +++ The header part of the document +++ -->
5  <front>
6  <to>Anna, Bernard, Didier, Johanna</to>
7  <date>Next Friday Evening at 8 pm</date>
8  <where>The Web Cafe</where>
9  <why>My first XML baby</why>
10 </front>
11 <!-- +++++ The main part of the document +++++ -->
12 <body>
13 <par>
14 I would like to invite you all to celebrate
15 the birth of <emph>Invitation</emph>, my
```

```

16 first XML document child.
17 </par>
18 <par>
19 Please do your best to come and join me next Friday
20 evening. And, do not forget to bring your friends.
21 </par>
22 <par>
23 I <emph>really</emph> look forward to see you soon!
24 </par>
25 </body>
26 <!-- +++ The closing part of the document +++ -->
27 <back>
28 <signature>Michel</signature>
29 </back>
30 </invitation>

```

La sintaxis del lenguaje queda reflejada en la siguiente *Document Type Definition (DTD)* que aparece en la sección 6.4.3 del libro de Goosens:

```

pl@nereida:~/Lp10910/Practicas/161009/src$ cat -n invitation.dtd
 1 <!-- invitation DTD -->
 2 <!-- May 26th 1998 mg -->
 3 <!ELEMENT invitation (front, body, back) >
 4 <!ELEMENT front      (to, date, where, why?) >
 5 <!ELEMENT date       (#PCDATA) >
 6 <!ELEMENT to         (#PCDATA) >
 7 <!ELEMENT where       (#PCDATA) >
 8 <!ELEMENT why         (#PCDATA) >
 9 <!ELEMENT body       (par+) >
10 <!ELEMENT par        (#PCDATA|emph)* >
11 <!ELEMENT emph       (#PCDATA) >
12 <!ELEMENT back       (signature) >
13 <!ELEMENT signature  (#PCDATA) >

```

El objetivo de esta práctica es escribir un programa Perl que usando las extensiones para expresiones regulares presentes en la versión 5.10 reconozca el lenguaje anterior.

Véase también:

- The LaTeX Web Companion
- Examples from The LaTeX Web Companion (véanse los subdirectorios correspondientes a los capítulos 6 y 7)

3.8. Análisis Sintáctico con Expresiones Regulares Perl

3.8.1. Introducción al Análisis Sintáctico con Expresiones Regulares

Como se ha comentado en la sección 3.2.5 Perl 5.10 permite el reconocimiento de expresiones definidas mediante gramáticas recursivas, siempre que estas puedan ser analizadas por un analizador recursivo descendente. Sin embargo, las expresiones regulares Perl 5.10 hace difícil construir una representación del árbol de análisis sintáctico abstracto. Además, la necesidad de explicitar en la regexp los blancos existentes entre los símbolos hace que la descripción sea menos robusta y menos legible.

Ejemplo: Traducción de expresiones aritméticas en infijo a postfijo

El siguiente ejemplo muestra una expresión regular que traduce expresiones de diferencias en infijo a postfijo.

Se usa una variable `$tran` para calcular la traducción de la subexpresión vista hasta el momento.

La gramática original que consideramos es recursiva a izquierdas:

```
exp ->  exp '-' digits
      |  digits
```

aplicando las técnicas explicadas en 4.8.1 y en el nodo de perlmonks 553889 transformamos la gramática en:

```
exp ->  digits rest
rest ->  '-' rest
      |  # empty
```

Sigue el código:

```
pl@nereida:~/Lperltesting$ cat -n infixtopostfix.pl
 1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
 2  use v5.10;
 3
 4  # Infix to postfix translator using 5.10 regexp
 5  # original grammar:
 6  #  exp ->  exp '-' digits
 7  #          | digits
 8  #
 9  # Applying left-recursion elimination we have:
10  #  exp ->  digits rest
11  #  rest ->  '-' rest
12  #          | # empty
13  #
14  my $input;
15  local our $tran = '';
16
17  my $regexp = qr{
18      (?&exp)
19
20      (?(DEFINE)
21          (?<exp>      ((?&digits)) \s* (?{ $tran .= "$^N "; say "tran=$tran"; }) (?&rest)
22                      (?{
23                          say "exp -> digits($^N) rest";
24                      })
25          )
26
27          (?<rest>      \s* - ((?&digits)) (?{ $tran .= "$^N - "; say "tran=$tran"; }) (?&
28                      (?{
29                          say "rest -> - digits($^N) rest";
30                      })
31                      |  # empty
32                      (?{
33                          say "rest -> empty";
34                      })
35          )
36  }
```

```

37         (?<digits> \s* (\d+)
38         )
39     )
40 }xms;
41
42 $input = <>;
43 chomp($input);
44 if ($input =~ $regexp) {
45     say "matches: $&\ntran=$tran";
46 }
47 else {
48     say "does not match";
49 }

```

La variable `$^N` contiene el valor que casó con el último paréntesis. Al ejecutar el código anterior obtenemos:

Véase la ejecución:

```

pl@nereida:~/Lperltesting$ ./infixtopostfix.pl
ab 5 - 3 -2 cd;
tran= 5
tran= 5 3 -
tran= 5 3 - 2 -
rest -> empty
rest -> - digits(2) rest
rest -> - digits( 3) rest
exp -> digits( 5) rest
matches: 5 - 3 -2
tran= 5 3 - 2 -

```

Como se ve, el recorrido primero profundo se traduce en la reconstrucción de una derivación a derechas.

Accediendo a los atributos de paréntesis anteriores mediante acciones intermedias

Es difícil extender el ejemplo anterior a lenguajes mas complejos debido a la limitación de que sólo se dispone de acceso al último paréntesis vía `$^N`. En muchos casos es necesario poder acceder a paréntesis/atributos anteriores.

El siguiente código considera el caso de expresiones con sumas, restas, multiplicaciones y divisiones. Utiliza la variable `op` y una acción intermedia (líneas 51-53) para almacenar el segundo paréntesis necesitado:

```

pl@nereida:~/Lperltesting$ cat -n ./calc510withactions3.pl
1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
2  use v5.10;
3
4  # Infix to postfix translator using 5.10 regexp
5  # Original grammar:
6
7  # exp ->  exp [-+] term
8  #       | term
9  # term -> term [*/] digits
10 #       | digits
11
12 # Applying left-recursion elimination we have:
13

```

```

14  # exp -> term re
15  # re ->  [+ -] term re
16  #          | # empty
17  # term -> digits rt
18  # rt ->  [* /] rt
19  #          | # empty
20
21
22  my $input;
23  my @stack;
24
25  local our $op = '';
26  my $regexp = qr{
27      (?&exp)
28
29      (? (DEFINE)
30          (?<exp>      (?&term) (?&re)
31                      (?{ say "exp -> term re" })
32          )
33
34          (?<re>      \s* ([+ -]) (?&term) \s* (?{ push @stack, $^N }) (?&re)
35                      (?{ say "re -> [+ -] term re" })
36          | # empty
37            (?{ say "re -> empty" })
38          )
39
40          (?<term>    ((?&digits))
41                      (?{ # intermediate action
42                          push @stack, $^N
43                      })
44          (?&rt)
45          (?{
46              say "term-> digits($^N) rt";
47          })
48          )
49
50          (?<rt>      \s* ([* /])
51                      (?{ # intermediate action
52                          local $op = $^N;
53                      })
54          ((?&digits)) \s*
55          (?{ # intermediate action
56              push @stack, $^N, $op
57          })
58          (?&rt) # end of <rt> definition
59          (?{
60              say "rt -> [* /] digits($^N) rt"
61          })
62          | # empty
63            (?{ say "rt -> empty" })
64          )
65
66          (?<digits> \s* \d+

```

```

67         )
68     )
69     }xms;
70
71     $input = <>;
72     chomp($input);
73     if ($input =~ $regexp) {
74         say "matches: $$\nStack=(@stack)";
75     }
76     else {
77         say "does not match";
78     }

```

Sigue una ejecución:

```

pl@nereida:~/Lperltesting$ ./calc510withactions3.pl
5-8/4/2-1
rt -> empty
term-> digits(5) rt
rt -> empty
rt -> [*/] digits(2) rt
rt -> [*/] digits(4) rt
term-> digits(8) rt
rt -> empty
term-> digits(1) rt
re -> empty
re -> [+~] term re
re -> [+~] term re
exp -> term re
matches: 5-8/4/2-1
Stack=(5 8 4 / 2 / - 1 -)

```

Accediendo a los atributos de paréntesis anteriores mediante @-

Sigue una solución alternativa que obvia la necesidad de introducir incómodas acciones intermedias. Utilizamos las variables @- y @+:

Since Perl 5.6.1 the special variables @- and @+ can functionally replace \$', \$\$ and \$'. These arrays contain pointers to the beginning and end of each match (see perlvar for the full story), so they give you essentially the same information, but without the risk of excessive string copying.

Véanse los párrafos en las páginas 89, 89) y 90 para mas información sobre @- y @+.

Nótese la función rc en las líneas 21-28. rc(1) nos retorna lo que casó con el último paréntesis, rc(2) lo que casó con el penúltimo, etc.

```

pl@nereida:~/Lperltesting$ cat -n calc510withactions4.pl
1    #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
2    use v5.10;
3
4    # Infix to postfix translator using 5.10 regexp
5    # Original grammar:
6
7    #  exp ->    exp [+~] term
8    #          | term
9    #  term ->  term [*/] digits

```

```

10      #          | digits
11
12      # Applying left-recursion elimination we have:
13
14      # exp -> term re
15      # re  ->  [+ -] term re
16      #          | # empty
17      # term -> digits rt
18      # rt   ->  [* /] rt
19      #          | # empty
20
21      sub rc {
22          my $ofs = - shift;
23
24          # Number of parenthesis that matched
25          my $np = @-;
26          #          string, ofsset, length
27          substr($_, $-[$ofs], $+[$np+$ofs] - $-[$ofs])
28      }
29
30      my $input;
31      my @stack;
32
33      my $regex = qr{
34          (?&exp)
35
36          (? (DEFINE)
37              (?<exp>      (?&term) (?&re)
38                          (?{ say "exp -> term re" })
39                      )
40
41              (?<re>      \s* ([+-]) (?&term) \s* (?{ push @stack, rc(1) }) (?&re)
42                          (?{ say "re -> [+ -] term re" })
43                          | # empty
44                          (?{ say "re -> empty" })
45                      )
46
47              (?<term>    ((?&digits))
48                          (?{ # intermediate action
49                              push @stack, rc(1)
50                          })
51                          (?&rt)
52                          (?{
53                              say "term-> digits(\".rc(1).\") rt";
54                          })
55                      )
56
57              (?<rt>      \s* ([*/]) ((?&digits)) \s*
58                          (?{ # intermediate action
59                              push @stack, rc(1), rc(2)
60                          })
61                          (?&rt) # end of <rt> definition
62                          (?{

```

```

63                                     say "rt -> [*/] digits(".rc(1).") rt"
64                                     })
65             | # empty
66             ({ say "rt -> empty" })
67         )
68
69         (?<digits> \s* \d+
70         )
71     )
72 }xms;
73
74 $input = <>;
75 chomp($input);
76 if ($input =~ $regex) {
77     say "matches: $$\nStack=(@stack)";
78 }
79 else {
80     say "does not match";
81 }

```

Ahora accedemos a los atributos asociados con los dos paréntesis, en la regla de <rt> usando la función rc:

```

(?<rt> \s*([*/]) ((?&digits)) \s*
      ({ # intermediate action
        push @stack, rc(1), rc(2)
      })

```

Sigue una ejecución del programa:

```

pl@nereida:~/Lperltesting$ ./calc510withactions4.pl
5-8/4/2-1
rt -> empty
term-> digits(5) rt
rt -> empty
rt -> [*/] digits(2) rt
rt -> [*/] digits(4) rt
term-> digits(8) rt
rt -> empty
term-> digits(1) rt
re -> empty
re -> [+ -] term re
re -> [+ -] term re
exp -> term re
matches: 5-8/4/2-1
Stack=(5 8 4 / 2 / - 1 -)
pl@nereida:~/Lperltesting$

```

Accediendo a los atributos de paréntesis anteriores mediante paréntesis con nombre

Una nueva solución: dar nombre a los paréntesis y acceder a los mismos:

```

47         (?<rt> \s*(?<op>[*/]) (?<num>(?!&digits)) \s*
48         ({ # intermediate action
49             push @stack, ${+{num}}, ${+{op}}
50         })

```

Sigue el código completo:

```
pl@nereida:~/Lperltesting$ cat -n ./calc510withnamedpar.pl
1  #!/usr/local/lib/perl/5.10.1/bin//perl5.10.1
2  use v5.10;
3
4  # Infix to postfix translator using 5.10 regexp
5  # Original grammar:
6
7  # exp ->  exp [-+] term
8  #       | term
9  # term ->  term [*/] digits
10 #        | digits
11
12 # Applying left-recursion elimination we have:
13
14 # exp ->  term re
15 # re  ->  [+ -] term re
16 #       | # empty
17 # term ->  digits rt
18 # rt  ->  [*/] rt
19 #       | # empty
20
21 my @stack;
22
23 my $regexp = qr{
24     (?&exp)
25
26     (? (DEFINE)
27         (?<exp>    (?&term) (?&re)
28                     (?{ say "exp -> term re" })
29                 )
30
31         (?<re>      \s* ([+-]) (?&term) \s* (?{ push @stack, $^N }) (?&re)
32                     (?{ say "re -> [+ -] term re" })
33                 | # empty
34                     (?{ say "re -> empty" })
35             )
36
37         (?<term>    ((?&digits))
38                     (?{ # intermediate action
39                         push @stack, $^N
40                     })
41                     (?&rt)
42                     (?{
43                         say "term-> digits($^N) rt";
44                     })
45             )
46
47         (?<rt>      \s*(?<op>[*/]) (?<num>(?!&digits)) \s*
48                     (?{ # intermediate action
49                         push @stack, ${num}, ${op}
50                     })
51                     (?&rt) # end of <rt> definition
```

```

52             (?{
53                 say "rt -> [*/] digits($^N) rt"
54             })
55         | # empty
56         (?{ say "rt -> empty" })
57     )
58
59     (?<digits> \s* \d+
60     )
61 )
62 }xms;
63
64 my $input = <>;
65 chomp($input);
66 if ($input =~ $regex) {
67     say "matches: $$\nStack=(@stack)";
68 }
69 else {
70     say "does not match";
71 }

```

Ejecución:

```

pl@nereida:~/Lperltesting$ ./calc510withnamedpar.pl
5-8/4/2-1
rt -> empty
term-> digits(5) rt
rt -> empty
rt -> [*/] digits(2) rt
rt -> [*/] digits(4) rt
term-> digits(8) rt
rt -> empty
term-> digits(1) rt
re -> empty
re -> [+ -] term re
re -> [+ -] term re
exp -> term re
matches: 5-8/4/2-1
Stack=(5 8 4 / 2 / - 1 -)

```

Véase También

- El nodo *Backreference variables in code embedded inside Perl 5.10 regexps* en PerlMonks
- El nodo *Strange behavior of @- and @+ in perl5.10 regexps* en PerlMonks

3.8.2. Construyendo el AST con Expresiones Regulares 5.10

Construiremos en esta sección un traductor de infijo a postfijo utilizando una aproximación general: construiremos una representación del Abstract Syntax Tree o AST (véase la sección 4.9 Árbol de Análisis Abstracto para una definición detallada de que es un árbol sintáctico).

Como la aplicación es un poco mas compleja la hemos dividido en varios ficheros. Esta es la estructura:

.


```

|-- ASTandtrans3.pl      # programa principal
|-- BinaryOp.pm         # clases para el manejo de los nodos del AST
|-- testreegxpparen.pl  # prueba para Regexp::Paren
'-- Regexp
    '-- Paren.pm        # módulo de extensión de $^N

```

La salida del programa puede ser dividida en tres partes. La primera muestra una antiderivación a derechas inversa:

```

pl@nereida:~/Lperltesting$ ./ASTandtrans3.pl
2*(3-4)
factor -> NUM(2)
factor -> NUM(3)
rt -> empty
term-> factor rt
factor -> NUM(4)
rt -> empty
term-> factor rt
re -> empty
re -> [+ -] term re
exp -> term re
factor -> ( exp )
rt -> empty
rt -> [*/] factor rt
term-> factor rt
re -> empty
exp -> term re
matches: 2*(3-4)

```

Que leída de abajo a arriba nos da una derivación a derechas de la cadena $2*(3-4)$:

```

exp => term re => term => factor rt =>
factor [*/](*) factor rt => factor [*/](*) factor =>
factor [*/](*) ( exp ) => factor [*/](*) ( term re ) =>
factor [*/](*) ( term [+ -](-) term re ) =>
factor [*/](*) ( term [+ -](-) term ) =>
factor [*/](*) ( term [+ -](-) factor rt ) =>
factor [*/](*) ( term [+ -](-) factor ) =>
factor [*/](*) ( term [+ -](-) NUM(4) ) =>
factor [*/](*) ( factor rt [+ -](-) NUM(4) ) =>
factor [*/](*) ( factor [+ -](-) NUM(4) ) =>
factor [*/](*) ( NUM(3) [+ -](-) NUM(4) ) =>
NUM(2) [*/](*) ( NUM(3) [+ -](-) NUM(4) )

```

La segunda parte nos muestra la representación del AST para la entrada dada $2*(3-4)$:

```

AST:
$VAR1 = bless( {
    'left' => bless( {
    'right' => bless( {
    'left' => bless( {
    'right' => bless( {
    'op' => '- '
    }, 'ADD' ),
    'op' => '*'
    }, 'MULT' );

```

La última parte de la salida nos muestra la traducción a postfijo de la expresión en infijo suministrada en la entrada (2*(3-4)):

```
2 3 4 - *
```

Programa Principal: usando la pila de atributos

La gramática original que consideramos es recursiva a izquierdas:

```
exp    ->  exp [-+] term
        | term
term   ->  term [*/] factor
        | factor
factor ->  \( exp \)
        | \(d+
```

aplicando las técnicas explicadas en 4.8.2 es posible transformar la gramática en una no recursiva por la izquierda:

```
exp      ->  term restoexp
restoexp ->  [-+] term restoexp
        | # vacío
term     ->  term restoterm
restoterm ->  [*/] factor restoterm
        | # vacío
factor   ->  \( exp \)
        | \(d+
```

Ahora bien, no basta con transformar la gramática en una equivalente. Lo que tenemos como punto de partida no es una gramática sino un *esquema de traducción* (véase la sección 4.7) que construye el AST asociado con la expresión. Nuestro esquema de traducción conceptual es algo así:

```
exp    ->  exp ([-+]) term      { ADD->new(left => $exp, right => $term, op => $1) }
        | term                  { $term }
term   ->  term ([*/]) factor   { MULT->new(left => $exp, right => $term, op => $1) }
        | factor                { $factor }
factor ->  \( exp \)            { $exp }
        | (\(d+)                { NUM->new(val => $1) }
```

Lo que queremos conseguir un conjunto de acciones semánticas asociadas para gramática no recursiva que sea equivalente a este.

Este es el programa resultante una vez aplicadas las transformaciones. La implementación de la asociación entre símbolos y atributos la realizamos manualmente mediante una pila de atributos:

```
pl@nereida:~/Lperltesting$ cat -n ./ASTandtrans3.pl
1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
2  use v5.10;
3  use strict;
4  use Regexp::Paren qw{g};
5  use BinaryOp;
6
7  use Data::Dumper;
8  $Data::Dumper::Indent = 1;
9
10 # Builds AST
11 my @stack;
12 my $regex = qr{
```

```

13      (?&exp)
14
15      (? (DEFINE)
16          (?<exp>      (?&term) (?&re)
17                      (?{ say "exp -> term re" })
18          )
19
20          (?<re>      \s* ([+-]) (?&term)
21                      (?{ # intermediate action
22                          local our ($ch1, $term) = splice @stack, -2;
23
24                          push @stack, ADD->new( {left => $ch1, right => $term, op => g(1)
25                          })
26                      (?&re)
27                      (?{ say "re -> [+ -] term re" })
28          | # empty
29                      (?{ say "re -> empty" })
30          )
31
32          (?<term>      ((?&factor)) (?&rt)
33                      (?{
34                          say "term-> factor rt";
35                      })
36          )
37
38          (?<rt>      \s* ([*/]) (?&factor)
39                      (?{ # intermediate action
40                          local our ($ch1, $ch2) = splice @stack, -2;
41
42                          push @stack, MULT->new({left => $ch1, right => $ch2, op => g(1)
43                          })
44                      (?&rt) # end of <rt> definition
45                      (?{
46                          say "rt -> [*/] factor rt"
47                      })
48          | # empty
49                      (?{ say "rt -> empty" })
50          )
51
52          (?<factor> \s* (\d+)
53                      (?{
54                          say "factor -> NUM($^N)";
55                          push @stack, bless { 'val' => g(1) }, 'NUM';
56                      })
57          | \s* \ ( (?&exp) \s* \ )
58                      (?{ say "factor -> ( exp )" })
59          )
60      )
61 }xms;
62
63 my $input = <>;
64 chomp($input);
65 if ($input =~ $regex) {

```

```

66     say "matches: $&";
67     my $ast = pop @stack;
68     say "AST:\n", Dumper $ast;
69
70     say $ast->translate;
71 }
72 else {
73     say "does not match";
74 }

```

Las Clases representando a los AST

Cada nodo del AST es un objeto. La clase del nodo nos dice que tipo de nodo es. Así los nodos de la clase **MULT** agrupan a los nodos de multiplicación y división. Los nodos de la clase **ADD** agrupan a los nodos de suma y resta. El procedimiento general es asociar un método **translate** con cada clase de nodo. De esta forma se logra el polimorfismo necesario: cada clase de nodo sabe como traducirse y el método **translate** de cada clase puede escribirse como

- Obtener los resultados de llamar a **\$child->translate** para cada uno de los nodos hijos **\$child**. Por ejemplo, si el nodo fuera un nodo **IF_ELSE** de un hipotético lenguaje de programación, se llamaría a los métodos **translate** sobre sus tres hijos **boolexpr**, **ifstatement** y **elstatement**.
- Combinar los resultados para producir la traducción adecuada del nodo actual.

Es esta combinación la que mas puede cambiar según el tipo de nodo. Así, en el caso de el nodo **IF_ELSE** el pseudocódigo para la traducción sería algo parecido a esto:

```

my $self = shift;
my $etiqueta1 = generar_nueva_etiqueta;
my $etiqueta2 = generar_nueva_etiqueta;

my $boolexpr      = $self->boolexpr->translate;
my $ifstatement   = $self->ifstatement->translate,
my $elstatement   = $self->elstatement->translate,
return << "ENDTRANS";
    $boolexpr
    JUMPZERO $etiqueta1:
    $ifstatement
    JUMP      $etiqueta2:
$etiqueta1:
    $elstatement
$etiqueta2:
ENDTRANS

```

Siguiendo estas observaciones el código de **BinaryOp.pm** queda así:

```

pl@nereida:~/Lperltesting$ cat -n BinaryOp.pm
1  package BinaryOp;
2  use strict;
3  use base qw(Class::Accessor);
4
5  BinaryOp->mk_accessors(qw{left right op});
6
7  sub translate {
8      my $self = shift;
9

```

```

10     return $self->left->translate." ".$self->right->translate." ".$self->op;
11 }
12
13 package ADD;
14 use base qw{BinaryOp};
15
16 package MULT;
17 use base qw{BinaryOp};
18
19 package NUM;
20
21 sub translate {
22     my $self = shift;
23
24     return $self->{val};
25 }
26
27 1;

```

Véase también:

- `Class::Accessor`

Accediendo a los paréntesis lejanos: El módulo `Regexp::Paren`

En esta solución utilizamos las variables `@-` y `@+` para construir una función que nos permite acceder a lo que pasó con los últimos paréntesis con memoria:

Since Perl 5.6.1 the special variables `@-` and `@+` can functionally replace `$'`, `$&` and `$'`. These arrays contain pointers to the beginning and end of each match (see `perlvar` for the full story), so they give you essentially the same information, but without the risk of excessive string copying.

Véanse los párrafos en las páginas 89, 89) y 90 para mas información sobre `@-` y `@+`.

`g(1)` nos retorna lo que pasó con el último paréntesis, `g(2)` lo que pasó con el penúltimo, etc.

```

pl@nereida:~/Lperltesting$ cat -n Regexp/Paren.pm
 1  package Regexp::Paren;
 2  use strict;
 3
 4  use base qw{Exporter};
 5
 6  our @EXPORT_OK = qw{g};
 7
 8  sub g {
 9      die "Error in 'Regexp::Paren::g'. Not used inside (?{ code }) construct\n" unless defined;
10      my $ofs = - shift;
11
12      # Number of parenthesis that matched
13      my $np = @-;
14      die "Error. Illegal 'Regexp::Paren::g' ref inside (?{ code }) construct\n" unless ($np > 0);
15      # $_ contains the string being matched
16      substr($_, $-[$ofs], $+[$np+$ofs] - $-[$ofs])
17  }
18
19  1;

```

```

20 =head1 NAME
21
22
23 Regexp::Paren - Extends $^N inside (?{ ... }) constructs
24
25 =head1 SYNOPSIS
26
27 use Regexp::Paren qw{g};
28
29 'abcde' =~ qr{(.) (.) (.)
30             (?{ print g(1)." ".g(2)." ".g(3)." \n" })           # c b a
31             (.)         (?{ print g(1)." ".g(2)." ".g(3)." ".g(4)." \n" }) # d c b
32             (.)         (?{ print g(1)." ".g(2)." ".g(3)." ".g(4)." ".g(5)." \n" }) # e d c
33             }x;
34
35 print g(1)." ".g(2)." ".g(3)." ".g(4)." ".g(5)." \n"; # error!
36
37 =head1 DESCRIPTION
38
39 Inside a C<(?...)> construct, C<g(1)> refers to what matched the last parenthesis
40 (like C<$^N>), C<g(2)> refers to the string that matched with the parenthesis before
41 the last, C<g(3)> refers to the string that matched with the parenthesis at distance 3,
42 etc.
43
44 =head1 SEE ALSO
45
46 =over 2
47
48 =item * L<perlre>
49
50 =item * L<perlretut>
51
52 =item * PerlMonks node I<Strange behavior o> C<@-> I<and> C<@+> I<in perl5.10 regexps> L<http://perlmonks.org/?q=51111>
53
54 =item * PerlMonks node I<Backreference variables in code embedded inside Perl 5.10 regexps> L<http://perlmonks.org/?q=51111>
55
56 =back
57
58 =head1 AUTHOR
59
60 Casiano Rodriguez-Leon (casiano@ull.es)
61
62 =head1 ACKNOWLEDGMENTS
63
64 This work has been supported by CEE (FEDER) and the Spanish Ministry of
65 I<Educacion y Ciencia> through I<Plan Nacional I+D+I> number TIN2005-08818-C04-04
66 (ULL::OPLINK project L<http://www.oplink.ull.es/>).
67 Support from Gobierno de Canarias was through GC02210601
68 (I<Grupos Consolidados>).
69 The University of La Laguna has also supported my work in many ways
70 and for many years.
71
72 =head1 LICENCE AND COPYRIGHT

```

73
74 Copyright (c) 2009- Casiano Rodriguez-Leon (casiano@cull.es). All rights reserved.
75
76 These modules are free software; you can redistribute it and/or
77 modify it under the same terms as Perl itself. See L<perlartistic>.
78
79 This program is distributed in the hope that it will be useful,
80 but WITHOUT ANY WARRANTY; without even the implied warranty of
81 MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.

Al ejecutar `perldoc Regexp::Paren` podemos ver la documentación incluida (véase la documentación en `perlpod` y `perlpodspec` así como la sección La Documentación en Perl para mas detalles):

NAME

Regexp::Paren - Extends $\N inside `(?{ ... })` constructs

SYNOPSIS

```
use Regexp::Paren qw{g};

'abcde' =~ qr{(.)(.)(.)
              (?{ print g(1)." ".g(2)." ".g(3)."\\n" })           # c
              (.)        (?{ print g(1)." ".g(2)." ".g(3)." ".g(4)."\\n" }) # d
              (.)        (?{ print g(1)." ".g(2)." ".g(3)." ".g(4)." ".g(5)."\\n" }) # e
              }x;

print g(1)." ".g(2)." ".g(3)." ".g(4)." ".g(5)."\\n"; # error!
```

DESCRIPTION

Inside a `(?{ ... })` construct, `g(1)` refers to what matched the last parenthesis (like $\N), `g(2)` refers to the string that matched with the parenthesis before the last, `g(3)` refers to the string that matched with the parenthesis at distance 3, etc.

SEE ALSO

- * perlre
- * perlretut
- * PerlMonks node *Strange behavior of "@-" *and* "@+" *in perl5.10 regexps** <http://www.perlmonks.org/?node_id=794736>
- * PerlMonks node *Backreference variables in code embedded inside Perl 5.10 regexps** <http://www.perlmonks.org/?node_id=794424>

AUTHOR

Casiano Rodriguez-Leon (casiano@ull.es)

ACKNOWLEDGMENTS

This work has been supported by CEE (FEDER) and the Spanish Ministry of *Educacion y Ciencia* through *Plan Nacional I+D+I* number TIN2005-08818-C04-04 (ULL::OPLINK project <<http://www.oplink.ull.es/>>). Support from Gobierno de Canarias was through GC02210601 (*Grupos Consolidados*). The University of La Laguna has also supported my work in many ways and for many years.

LICENCE AND COPYRIGHT

Copyright (c) 2009- Casiano Rodriguez-Leon (casiano@ull.es). All rights

3.9. Práctica: Traducción de invitation a HTML

Esta práctica es continuación de la práctica *un lenguaje para componer invitaciones* especificada en la sección 3.7.

El objetivo es traducir la entrada escrita en el lenguaje de invitaciones a HTML. La traducción del

ejemplo anterior debería ser parecida a esta:

```
pl@nereida:~/Lp10910/Practicas/161009/src$ cat -n invitati
1  <?xml version="1.0"?>
2  <!DOCTYPE invitation SYSTEM "invitation.dtd">
3  <invitation>
4  <!-- ++++ The header part of the document ++++ -->
5  <front>
6  <to>Anna, Bernard, Didier, Johanna</to>
7  <date>Next Friday Evening at 8 pm</date>
8  <where>The Web Cafe</where>
9  <why>My first XML baby</why>
10 </front>
11 <!-- +++++ The main part of the document +++++ -->
12 <body>
13 <par>
14 I would like to invite you all to celebrate
15 the birth of <emph>Invitation</emph>, my
16 first XML document child.
17 </par>
18 <par>
19 Please do your best to come and join me next Friday
20 evening. And, do not forget to bring your friends.
21 </par>
22 <par>
23 I <emph>really</emph> look forward to see you soon!
24 </par>
25 </body>
26 <!-- +++ The closing part of the document +++ -->
27 <back>
28 <signature>Michel</signature>
29 </back>
30 </invitation>
```

Para ver el resultado en su navegador visite el fichero invitation.html

Su programa deberá producir un Abstract Syntax Tree. Los nodos serán objetos. Cada clase (FRONT, TO, etc.) deberá de disponer de un método `translate`.

Para simplificar el proceso de traducción a HTML se sugiere utilizar una hoja de estilo parecida a la siguiente (tomada de la sección 7.4.4 del citado libro de Goosens):

```
pl@nereida:~/Lp10910/Practicas/161009/src$ cat -n invit.css
1  /* CSS stylesheet for invitation1 in HTML */
2  BODY {margin-top: 1em;      /* global page parameters */
3      margin-bottom: 1em;
4      margin-left: 1em;
5      margin-right: 1em;
6      font-family: serif;
7      line-height: 1.1;
8      color: black;
9  }
```

```

10 H1    {text-align: center; /* for global title */
11        font-size: x-large;
12    }
13 P      {text-align: justify; /* paragraphs in body */
14        margin-top: 1em;
15    }
16 TABLE { border-width: 0pt }
17 TBODY  { border-width: 0pt }
18 TD[class="front"] { /* table data in front matter */
19        text-align: left;
20        font-weight: bold;
21    }
22 TD.front { /* table data in front matter */
23        text-align: left;
24        font-weight: bold;
25    }
26 EM     {font-style: italic; /* emphasis in body */
27    }
28 P.signature { /* signature */
29        text-align: right;
30        font-weight: bold;
31    }

```

Véase también:

- The LaTeX Web Companion
- Examples from The LaTeX Web Companion (véanse los subdirectorios correspondientes a los capítulos 6 y 7)
- CSS Tutorial
- Edición extremadamente simple de HTML
- Perl-XML Frequently Asked Questions

3.10. Análisis Sintáctico con Regexp::Grammars

El módulo Regexp::Grammars escrito por Damian Conway extiende las expresiones regulares Perl con la capacidad de generar representaciones del árbol de análisis sintáctico abstracto y obviando la necesidad de explicitar los blancos. El módulo necesita para funcionar una versión de Perl superior o igual a la 5.10.

3.10.1. Introducción

El Problema

La documentación de Regexp::Grammars establece cual es el problema que aborda el módulo:

...Perl5.10 makes possible to use regexes to recognize complex, hierarchical—and even recursive—textual structures. The problem is that Perl 5.10 doesn't provide any support for extracting that hierarchical data into nested data structures. In other words, using Perl 5.10 you can match complex data, but not parse it into an internally useful form.

An additional problem when using Perl 5.10 regexes to match complex data formats is that you have to make sure you remember to insert whitespace- matching constructs (such as \s) at every possible position where the data might contain ignorable whitespace. This reduces the readability of such patterns, and increases the chance of errors (typically caused by overlooking a location where whitespace might appear).*

Una solución: Regexp::Grammars

The Regexp::Grammars module solves both those problems.

If you import the module into a particular lexical scope, it preprocesses any regex in that scope, so as to implement a number of extensions to the standard Perl 5.10 regex syntax. These extensions simplify the task of defining and calling subrules within a grammar, and allow those subrule calls to capture and retain the components of they match in a proper hierarchical manner.

La sintaxis de una expresión regular Regexp::Grammars

Las expresiones regulares Regexp::Grammars aumentan las regex Perl 5.10. La sintaxis se expande y se modifica:

A Regexp::Grammars specification consists of a pattern (which may include both standard Perl 5.10 regex syntax, as well as special Regexp::Grammars directives), followed by one or more rule or token definitions.

Sigue un ejemplo:

```
pl@nereida:~/Lregexpg grammars/demo$ cat -n balanced_brackets.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8     qr{
 9         (<pp>)
10
11         <rule: pp>    \( (? : [^()]*)+ | <escape> | <pp> )* \)
12
13         <token: escape> \\.
14
15     }xs;
16 };
17
18 while (my $input = <>) {
19     while ($input =~ m{$rbb}g) {
20         say("matches: <$&>");
21         say Dumper \% /;
22     }
23 }
```

Note that there is no need to explicitly place \s subpatterns throughout the rules; that is taken care of automatically.*

...

The initial pattern ((<pp>)) acts like the top rule of the grammar, and must be matched completely for the grammar to match.

The rules and tokens are declarations only and they are not directly matched. Instead, they act like subroutines, and are invoked by name from the initial pattern (or from within a rule or token).

Each rule or token extends from the directive that introduces it up to either the next rule or token directive, or (in the case of the final rule or token) to the end of the grammar.

El hash %/: Una representación del AST Al ejecutar el programa anterior con entrada $(2*(3+5))*4+(2-3)$ produce:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 balanced_brackets.pl
(2*(3+5))*4+(2-3)
matches: <(2*(3+5))>
$VAR1 = {
    '' => '(2*(3+5))',
    'pp' => {
        '' => '(2*(3+5))',
        'pp' => '(3+5)'
    }
};

matches: <(2-3)>
$VAR1 = {
    '' => '(2-3)',
    'pp' => '(2-3)'
};
```

Each rule calls the subrules specified within it, and then return a hash containing whatever result each of those subrules returned, with each result indexed by the subrule's name.

In this way, each level of the hierarchical regex can generate hashes recording everything its own subrules matched, so when the entire pattern matches, it produces a tree of nested hashes that represent the structured data the pattern matched.

...

In addition each result-hash has one extra key: the empty string. The value for this key is whatever string the entire subrule call matched.

Diferencias entre token y rule

The difference between a token and a rule is that a token treats any whitespace within it exactly as a normal Perl regular expression would. That is, a sequence of whitespace in a token is ignored if the /x modifier is in effect, or else matches the same literal sequence of whitespace characters (if /x is not in effect).

En el ejemplo anterior el comportamiento es el mismo si se reescribe la regla para el token `escape` como:

```
13      <rule: escape> \\\.
```

En este otro ejemplo mostramos que la diferencia entre token y rule es significativa:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n tokenvsrule.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8     qr{
 9         <s>
10
11         <rule: s> <a> <c>
```

```

12
13     <rule: c>  c d
14
15     <token: a>  a b
16
17     }xs;
18 };
19
20 while (my $input = <>) {
21     if ($input =~ m{$rbb}) {
22         say("matches: <$&>");
23         say Dumper \% /;
24     }
25     else {
26         say "Does not match";
27     }
28 }

```

Al ejecutar este programa vemos la diferencia en la interpretación de los blancos:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 tokenvsrule.pl
```

```
ab c d
```

```
matches: <ab c d>
```

```
$VAR1 = {
    '' => 'ab c d',
    's' => {
        '' => 'ab c d',
        'c' => 'c d',
        'a' => 'ab'
    }
};
```

```
a b c d
```

```
Does not match
```

```
ab cd
```

```
matches: <ab cd>
```

```
$VAR1 = {
    '' => 'ab cd',
    's' => {
        '' => 'ab cd',
        'c' => 'cd',
        'a' => 'ab'
    }
};
```

Obsérvese como la entrada a b c d es rechazada mientras que la entrada ab c d es aceptada.

Redefinición de los espacios en blanco

In a rule, any sequence of whitespace (except those at the very start and the very end of the rule) is treated as matching the implicit subrule <.ws>, which is automatically predefined to match optional whitespace (i.e. \s).*

You can explicitly define a <ws> token to change that default behaviour. For example, you could alter the definition of whitespace to include Perlsh comments, by adding an explicit <token: ws>:

```
<token: ws>
(?: \s+ | #[^\n]* )*
```

But be careful not to define <ws> as a rule, as this will lead to all kinds of infinitely recursive unpleasantness.

El siguiente ejemplo ilustra como redefinir <ws>:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n tokenvsruleandws.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8     no warnings 'uninitialized';
 9     qr{
10         <s>
11
12         <token: ws> (?: \s+ | /\* .*? \*/)*+
13
14         <rule: s> <a> <c>
15
16         <rule: c>  c d
17
18         <token: a>  a b
19
20     }xs;
21 };
22
23 while (my $input = <>) {
24     if ($input =~ m{$rbb}) {
25         say Dumper \% /;
26     }
27     else {
28         say "Does not match";
29     }
30 }
```

Ahora podemos introducir comentarios en la entrada:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 -w tokenvsruleandws.pl
ab /* 1 */ c d
$VAR1 = {
    '' => 'ab /* 1 */ c d',
    's' => {
        '' => 'ab /* 1 */ c d',
        'c' => 'c d',
        'a' => 'ab'
    }
};
```

Llamando a las subreglas

To invoke a rule to match at any point, just enclose the rule's name in angle brackets (like in Perl 6). There must be no space between the opening bracket and the rulename. For example:

```
qr{
    file:          # Match literal sequence 'f' 'i' 'l' 'e' ':'
    <name>          # Call <rule: name>
    <options>?      # Call <rule: options> (it's okay if it fails)

    <rule: name>
        # etc.
}x;
```

If you need to match a literal pattern that would otherwise look like a subrule call, just backslash-escape the leading angle:

```
qr{
    file:          # Match literal sequence 'f' 'i' 'l' 'e' ':'
    \<name>        # Match literal sequence '<' 'n' 'a' 'm' 'e' '>'
    <options>?      # Call <rule: options> (it's okay if it fails)

    <rule: name>
        # etc.
}x;
```

El siguiente programa ilustra algunos puntos discutidos en la cita anterior:

```
casiano@millo:~/src/perl/regexp-grammar-examples$ cat -n badbracket.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8     qr{
 9         (<pp>)
10
11         <rule: pp>    \( (?<b> > | \< | < escape> | <pp> )* \)
12
13         <token: b> > b
14
15         <token: escape> \\.
16
17     }xs;
18 };
19
20 while (my $input = <>) {
21     while ($input =~ m{$rbb}g) {
22         say("matches: <${&}>");
23         say Dumper \%/%;
24     }
25 }
```

Obsérvense los blancos en < escape> y en <token: b > b. Pese a ello el programa funciona:

```
casiano@millo:~/src/perl/regexp-grammar-examples$ perl5.10.1 badbracket.pl  
(\(\))
```

```
matches: <(\(\))>
```

```
$VAR1 = {  
    '' => '(\(\(\(\)\))',  
    'pp' => {  
        '' => '(\(\(\(\(\)\))',  
        'escape' => '\\\\'  
    }  
};
```

(b)

```
matches: <(b)>
```

```
$VAR1 = {  
    '' => '(b)',  
    'pp' => {  
        '' => '(b)',  
        'b' => 'b'  
    }  
};
```

(<)

```
matches: <(<)>
```

```
$VAR1 = {  
    '' => '(<)',  
    'pp' => '(<)',  
};
```

(c)

```
casiano@millo:
```

Eliminación del anidamiento de ramas unarias en%/

... Note, however, that if the result-hash at any level contains only the empty-string key (i.e. the subrule did not call any sub-subrules or save any of their nested result-hashes), then the hash is unpacked and just the matched substring itself is returned.

For example, if <rule: sentence> had been defined:

```
<rule: sentence>  
    I see dead people
```

then a successful call to the rule would only add:

```
sentence => 'I see dead people'
```

to the current result-hash.

This is a useful feature because it prevents a series of nested subrule calls from producing very unwieldy data structures. For example, without this automatic unpacking, even the simple earlier example:

```
<rule: sentence>  
    <noun> <verb> <object>
```


would produce something needlessly complex, such as:

```
sentence => {
  ""      => 'I saw a dog',
  noun    => {
    ""    => 'I',
  },
  verb    => {
    ""    => 'saw',
  },
  object  => {
    ""      => 'a dog',
    article => {
      ""    => 'a',
    },
    noun    => {
      ""    => 'dog',
    },
  },
}
```

El siguiente ejemplo ilustra este punto:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n unaryproductions.pl
1  use strict;
2  use warnings;
3  use 5.010;
4  use Data::Dumper;
5
6  my $rbb = do {
7      use Regexp::Grammars;
8      qr{
9          <s>
10
11          <rule: s> <noun> <verb> <object>
12
13          <token: noun> he | she | Peter | Jane
14
15          <token: verb> saw | sees
16
17          <token: object> a\s+dog | a\s+cat
18
19      }x;
20 };
21
22 while (my $input = <>) {
23     while ($input =~ m{$rbb}g) {
24         say("matches: <$&>");
25         say Dumper \% /;
26     }
27 }
```

Sigue una ejecución del programa anterior:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 unaryproductions.pl
```

```

he saw a dog
matches: <he saw a dog>
$VAR1 = {
    '' => 'he saw a dog',
    's' => {
        '' => 'he saw a dog',
        'object' => 'a dog',
        'verb' => 'saw',
        'noun' => 'he'
    }
};

Jane sees a cat
matches: <Jane sees a cat>
$VAR1 = {
    '' => 'Jane sees a cat',
    's' => {
        '' => 'Jane sees a cat',
        'object' => 'a cat',
        'verb' => 'sees',
        'noun' => 'Jane'
    }
};

```

Ámbito de uso de Regexp::Grammars

Cuando se usa Regexp::Grammars como parte de un programa que utiliza otras regexes hay que evitar que Regexp::Grammars procese las mismas. Regexp::Grammars reescribe las expresiones regulares durante la fase de preproceso. Esta por ello presenta las mismas limitaciones que cualquier otra forma de 'source filtering' (véase `perlfilter`). Por ello es una buena idea declarar la gramática en un bloque `do` restringiendo de esta forma el ámbito de acción del módulo.

```

5 my $calculator = do{
6     use Regexp::Grammars;
7     qr{
8         .....
9     }xms
10 };

```

3.10.2. Objetos

When a grammar has parsed successfully, the `%/` variable will contain a series of nested hashes (and possibly arrays) representing the hierarchical structure of the parsed data.

Typically, the next step is to walk that tree, extracting or converting or otherwise processing that information. If the tree has nodes of many different types, it can be difficult to build a recursive subroutine that can navigate it easily.

A much cleaner solution is possible if the nodes of the tree are proper objects. In that case, you just define a `translate()` method for each of the classes, and have every node call that method on each of its children. The chain of `translate()` calls would cascade down the nodes of the tree, each one invoking the appropriate `translate()` method according to the type of node encountered.

The only problem is that, by default, Regexp::Grammars returns a tree of plain-old hashes, not `Class::Whatever` objects. Fortunately, it's easy to request that the result has-

hes be automatically blessed into the appropriate classes, using the `<objrule:...>` and `<objtoken:...>` directives.

These directives are identical to the `<rule:...>` and `<token:...>` directives (respectively), except that the rule or token they create will also bless the hash it normally returns, converting it to an object of a class whose name is the same as the rule or token itself.

For example:

```
<objrule: Element>
# ...Defines a rule that can be called as <Element>
# ...and which returns a hash-based Element object
```

The IDENTIFIER of the rule or token may also be fully qualified. In such cases, the rule or token is defined using only the final short name, but the result object is blessed using the fully qualified long name. For example:

```
<objrule: LaTeX::Element>
# ...Defines a rule that can be called as <Element>
# ...and which returns a hash-based LaTeX::Element object
```

This can be useful to ensure that returned objects don't collide with other namespaces in your program.

Note that you can freely mix object-returning and plain-old-hash-returning rules and tokens within a single grammar, though you have to be careful not to subsequently try to call a method on any of the unblessed nodes.

3.10.3. Renombrando los resultados de una subregla

Nombre de la regla versus Nombre del Resultado

No siempre el nombre de la regla es el mas apropiado para ser el nombre del resultado:

It is not always convenient to have subrule results stored under the same name as the rule itself. Rule names should be optimized for understanding the behaviour of the parser, whereas result names should be optimized for understanding the structure of the data. Often those two goals are identical, but not always; sometimes rule names need to describe what the data looks like, while result names need to describe what the data means.

Colisión de nombres de reglas

For example, sometimes you need to call the same rule twice, to match two syntactically identical components whose positions give them semantically distinct meanings:

```
<rule: copy_cmd>
  copy <file> <file>
```

The problem here is that, if the second call to <file> succeeds, its result-hash will be stored under the key file, clobbering the data that was returned from the first call to <file>.

Aliasing

To avoid such problems, Regexp::Grammars allows you to alias any subrule call, so that it is still invoked by the original name, but its result-hash is stored under a different key. The syntax for that is: <alias=rulename>. For example:

```
<rule: copy_cmd>
  copy <from=file> <to=file>
```

Here, `<rule: file>` is called twice, with the first result-hash being stored under the key `from`, and the second result-hash being stored under the key `to`.

Note, however, that the alias before the `=` must be a proper identifier (i.e. a letter or underscore, followed by letters, digits, and/or underscores). Aliases that start with an underscore and aliases named `MATCH` have special meaning.

Normalización de los resultados mediante aliasing

Aliases can also be useful for normalizing data that may appear in different formats and sequences. For example:

```
<rule: copy_cmd>
  copy <from=file>      <to=file>
| dup   <to=file> as <from=file>
|       <from=file> ->  <to=file>
|       <to=file> <-  <from=file>
```

Here, regardless of which order the old and new files are specified, the result-hash always gets:

```
copy_cmd => {
  from => 'oldfile',
  to   => 'newfile',
}
```

Ejemplo

El siguiente programa ilustra los comentarios de la documentación:

```
pl@nereida:~/Lregexgrammars/demo$ cat -n copygrammar.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7   use Regexp::Grammars;
 8   qr{
 9     <copy_cmd>
10
11     <rule: copy_cmd>
12       copy <from=file> <to=file>
13       |   <from=file> ->  <to=file>
14       |   <to=file>  <-  <from=file>
15
16       <token: file> [\w./\\]+
17     }x;
18 };
19
20 while (my $input = <>) {
21   while ($input =~ m{$rbb}g) {
22     say("matches: <$$>");
23     say Dumper \%/%;
24   }
25 }
```

Cuando lo ejecutamos obtenemos:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 copygrammar.pl
copy a b
matches: <copy a b>
$VAR1 = {
    '' => 'copy a b',
    'copy_cmd' => {
        '' => 'copy a b',
        'to' => 'b',
        'from' => 'a'
    }
};

b <- a
matches: <b <- a>
$VAR1 = {
    '' => 'b <- a',
    'copy_cmd' => {
        '' => 'b <- a',
        'to' => 'b',
        'from' => 'a'
    }
};

a -> b
matches: <a -> b>
$VAR1 = {
    '' => 'a -> b',
    'copy_cmd' => {
        '' => 'a -> b',
        'to' => 'b',
        'from' => 'a'
    }
};
```

3.10.4. Listas

El operador de cierre positivo

If a subrule call is quantified with a repetition specifier:

```
<rule: file_sequence>
  <file>+
```

then each repeated match overwrites the corresponding entry in the surrounding rule's result-hash, so only the result of the final repetition will be retained. That is, if the above example matched the string foo.pl bar.py baz.php, then the result-hash would contain:

```
file_sequence {
    "" => 'foo.pl bar.py baz.php',
    file => 'baz.php',
}
```

Operadores de listas y espacios en blanco

Existe un caveat con el uso de los operadores de repetición y el manejo de los blancos. Véase el siguiente programa:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n numbers3.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8
 9     qr{
10         <numbers>
11
12         <rule: numbers>
13             (<number>)+
14
15         <token: number> \s*\d+
16     }xms;
17 };
18
19 while (my $input = <>) {
20     if ($input =~ m{$rbb}) {
21         say("matches: <$>");
22         say Dumper \% /;
23     }
24 }
```

Obsérvese el uso explícito de espacios `\s*\d+` en la definición de `number`.

Sigue un ejemplo de ejecución:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5_10_1 numbers3.pl
1 2 3 4
matches: <1 2 3 4>
$VAR1 = {
    '' => '1 2 3 4',
    'numbers' => {
        '' => '1 2 3 4',
        'number' => ' 4'
    }
};
```

Si se eliminan los blancos de la definición de `number`:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n numbers.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8
```

```

9      qr{
10         <numbers>
11
12         <rule: numbers>
13         (<number>)+
14
15         <token: number> \d+
16     }xms;
17 };
18
19 while (my $input = <>) {
20     if ($input =~ m{$rbb}) {
21         say("matches: <$&>");
22         say Dumper \% /;
23     }
24 }

```

se obtiene una conducta que puede sorprender:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 numbers.pl
12 34 56
matches: <12>
$VAR1 = {
    '' => '12',
    'numbers' => {
        '' => '12',
        'number' => '12'
    }
};

```

La explicación está en la documentación: véase la sección Grammar Syntax:

```

<rule: IDENTIFIER>
    Define a rule whose name is specified by the supplied identifier.
    Everything following the <rule:...> directive (up to the next <rule:...> or <token:...>
    directive) is treated as part of the rule being defined.
    Any whitespace in the rule is replaced by a call to the <.ws> subrule (which defaults
    to matching \s*, but may be explicitly redefined).

```

También podríamos haber resuelto el problema introduciendo un blanco explícito dentro del cierre positivo:

```

<rule: numbers>
    (<number> )+

<token: number> \d+

```

Una Solución al problema de recordar los resultados de una lista: El uso de brackets

Usually, that's not the desired outcome, so Regexp::Grammars provides another mechanism by which to call a subrule; one that saves all repetitions of its results.

A regular subrule call consists of the rule's name surrounded by angle brackets. If, instead, you surround the rule's name with <[...]> (angle and square brackets) like so:

```

<rule: file_sequence>
    <[file]>+

```

then the rule is invoked in exactly the same way, but the result of that submatch is pushed onto an array nested inside the appropriate result-hash entry. In other words, if the above example matched the same foo.pl bar.py baz.php string, the result-hash would contain:

```
file_sequence {
    "" => 'foo.pl bar.py baz.php',
    file => [ 'foo.pl', 'bar.py', 'baz.php' ],
}
```

Teniendo en cuenta lo dicho anteriormente sobre los blancos dentro de los cuantificadores, es necesario introducir blancos dentro del operador de repetición:

```
pl@nereida:~/Lregexpg grammars/demo$ cat -n numbers4.pl
```

```
1 use strict;
2 use warnings;
3 use 5.010;
4 use Data::Dumper;
5
6 my $rbb = do {
7     use Regexp::Grammars;
8
9     qr{
10         <numbers>
11
12         <rule: numbers>
13         (?: <[number]> )+
14
15         <token: number> \d+
16     }xms;
17 };
18
19 while (my $input = <>) {
20     if ($input =~ m{$rbb}) {
21         say("matches: <$>");
22         say Dumper \% /;
23     }
24 }
```

Al ejecutar este programa obtenemos:

```
pl@nereida:~/Lregexpg grammars/demo$ perl5_10_1 numbers4.pl
```

```
1 2 3 4
matches: <1 2 3 4
>
$VAR1 = {
    '' => '1 2 3 4
',
    'numbers' => {
        '' => '1 2 3 4
',
        'number' => [ '1', '2', '3', '4' ]
    }
};
```


Otra forma de resolver las colisiones de nombres: salvarlos en una lista

This listifying subrule call can also be useful for non-repeated subrule calls, if the same subrule is invoked in several places in a grammar. For example if a cmdline option could be given either one or two values, you might parse it:

```
<rule: size_option>
  -size <[size]> (?: x <[size]> )?
```

The result-hash entry for size would then always contain an array, with either one or two elements, depending on the input being parsed.

Sigue un ejemplo:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n sizes.pl
```

```
1  use strict;
2  use warnings;
3  use 5.010;
4  use Data::Dumper;
5
6  my $rbb = do {
7      use Regexp::Grammars;
8
9      qr{
10         <command>
11
12         <rule: command> ls <size_option>
13
14         <rule: size_option>
15             -size <[size]> (?: x <[size]> )?
16
17         <token: size> \d+
18     }x;
19 };
20
21 while (my $input = <>) {
22     while ($input =~ m{$rbb}g) {
23         say("matches: <$&>");
24         say Dumper \% /;
25     }
26 }
```

Veamos su comportamiento con diferentes entradas:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 sizes.pl
```

```
ls -size 4
matches: <ls -size 4
>
$VAR1 = {
    '' => 'ls -size 4'
},
    'command' => {
        'size_option' => {
            '' => '-size 4'
        }
    }
},
```

```

        'size' => [ '4' ]
      },
      '' => 'ls -size 4
    ,
      }
    };

ls -size 2x8
matches: <ls -size 2x8
>
$VAR1 = {
    '' => 'ls -size 2x8
  ,
    'command' => {
        'size_option' => {
            '' => '-size 2x8
        ,
            'size' => [ '2', '8' ]
        },
        '' => 'ls -size 2x8
    ,
    }
};

```

Aliasing de listas

Listifying subrules can also be given aliases, just like ordinary subrules. The alias is always specified inside the square brackets:

```

<rule: size_option>
  -size <[size=pos_integer]> (?: x <[size=pos_integer]> )?

```

Here, the sizes are parsed using the `pos_integer` rule, but saved in the result-hash in an array under the key `size`.

Sigue un ejemplo:

```

pl@nereida:~/Lregexpgrammars/demo$ cat -n aliasedsizes.pl
1  use strict;
2  use warnings;
3  use 5.010;
4  use Data::Dumper;
5
6  my $rbb = do {
7      use Regexp::Grammars;
8
9      qr{
10         <command>
11
12         <rule: command> ls <size_option>
13
14         <rule: size_option>
15             -size <[size=int]> (?: x <[size=int]> )?
16

```

```

17     <token: int> \d+
18   }x;
19 };
20
21 while (my $input = <>) {
22   while ($input =~ m{$rbb}g) {
23     say("matches: <$&>");
24     say Dumper \%/%;
25   }
26 }

```

Veamos el resultado de una ejecución:

```

pl@nereida:~/Lregexpgrammars/demo$ perl5.10.1 aliasedsizes.pl
ls -size 2x4
matches: <ls -size 2x4
>
$VAR1 = {
    '' => 'ls -size 2x4
',
    'command' => {
        'size_option' => {
            '' => '-size 2x4
',
            'size' => [
                '2',
                '4'
            ]
        },
        '' => 'ls -size 2x4
',
    },
};

```

Caveat: Cierres y Warnings

En este ejemplo aparece <number>+ sin corchetes ni paréntesis:

```

pl@nereida:~/Lregexpgrammars/demo$ cat -n numbers5.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7   use Regexp::Grammars;
 8
 9   qr{
10     <numbers>
11
12     <rule: numbers>
13     <number>+
14
15     <token: number> \d+
16   }xms;

```

```

17  };
18
19  while (my $input = <>) {
20      if ($input =~ m{$rbb}) {
21          say("matches: <$&>");
22          say Dumper \% /;
23      }
24  }

```

Este programa produce un mensaje de advertencia:

```

pl@nereida:~/Lregexpg grammars/demo$ perl5.10.1 numbers5.pl
warn | Repeated subrule <number>+ will only capture its final match
      | (Did you mean <[number]>+ instead?)
      |

```

Si se quiere evitar el mensaje y se está dispuesto a asumir la pérdida de los valores asociados con los elementos de la lista se deberán poner el operando entre paréntesis (con o sin memoria).

Esto es lo que dice la documentación sobre este warning:

Repeated subrule <rule> will only capture its final match
You specified a subrule call with a repetition qualifier, such as:

`<ListElem>*`

or:

`<ListElem>+`

Because each subrule call saves its result in a hash entry of the same name, each repeated match will overwrite the previous ones, so only the last match will ultimately be saved. If you want to save all the matches, you need to tell Regexp::Grammars to save the sequence of results as a nested array within the hash entry, like so:

`<[ListElem]>*`

or:

`<[ListElem]>+`

If you really did intend to throw away every result but the final one, you can silence the warning by placing the subrule call inside any kind of parentheses. For example:

`(<ListElem>)*`

or:

`(?: <ListElem>)+`

3.10.5. Pseudo sub-reglas

Subpatrones

Aliases can also be given to standard Perl subpatterns, as well as to code blocks within a regex. The syntax for subpatterns is:

`<ALIAS= (SUBPATTERN) >`

In other words, the syntax is exactly like an aliased subrule call, except that the rule name is replaced with a set of parentheses containing the subpattern. Any parentheses-capturing or non-capturing-will do.

The effect of aliasing a standard subpattern is to cause whatever that subpattern matches to be saved in the result-hash, using the alias as its key. For example:

`<rule: file_command>`

`<cmd=(mv|cp|ln)> <from=file> <to=file>`

Here, the `<cmd=(mv|cp|ln)>` is treated exactly like a regular `(mv|cp|ln)`, but whatever substring it matches is saved in the result-hash under the key 'cmd'.

Sigue un ejemplo:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n subpattern.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8
 9     qr{
10         <file_command>
11
12         <rule: file_command>
13
14         <cmd=(mv|cp|ln)> <from=([\w./]+)> <to=([\w./]+)>
15
16     }x;
17 };
18
19 while (my $input = <>) {
20     while ($input =~ m{$rbb}g) {
21         say("matches: <${&}>");
22         say Dumper \% /;
23     }
24 }
```

y una ejecución:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 subpattern.pl
mv a b
matches: <mv a b>
$VAR1 = {
```

```

'' => 'mv a b',
'file_command' => {
    '' => 'mv a b',
    'to' => 'b',
    'cmd' => 'mv',
    'from' => 'a'
}
};

cp c d
matches: <cp c d>
$VAR1 = {
    '' => 'cp c d',
    'file_command' => {
        '' => 'cp c d',
        'to' => 'd',
        'cmd' => 'cp',
        'from' => 'c'
    }
}

```

Bloques de código

The syntax for aliasing code blocks is:

```
<ALIAS= (?{ your($code->here) }) >
```

Note, however, that the code block must be specified in the standard Perl 5.10 regex notation: (?{...}). A common mistake is to write:

```
<ALIAS= { your($code->here) } >
```

instead, which will attempt to interpolate \$code before the regex is even compiled, as such variables are only protected from interpolation inside a (?{...}).

When correctly specified, this construct executes the code in the block and saves the result of that execution in the result-hash, using the alias as its key. Aliased code blocks are useful for adding semantic information based on which branch of a rule is executed. For example, consider the copy_cmd alternatives shown earlier:

```

<rule: copy_cmd>
    copy <from=file>          <to=file>
| dup    <to=file> as <from=file>
|        <from=file> ->    <to=file>
|        <to=file> <-    <from=file>

```

Using aliased code blocks, you could add an extra field to the result- hash to describe which form of the command was detected, like so:

```

<rule: copy_cmd>
    copy <from=file>          <to=file> <type=(?{ 'std' })>
| dup    <to=file> as <from=file> <type=(?{ 'rev' })>
|        <from=file> ->    <to=file> <type=(?{ 'fwd' })>
|        <to=file> <-    <from=file> <type=(?{ 'bwd' })>

```

Now, if the rule matched, the result-hash would contain something like:

```

copy_cmd => {
    from => 'oldfile',
    to => 'newfile',
    type => 'fwd',
}

```

El siguiente ejemplo ilustra lo dicho en la documentación. En la línea 15 hemos introducido una regla para el control de errores⁸:

```

pl@nereida:~/Lregexpggrammars/demo$ cat -n aliasedcodeblock2.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5
 6 my $rbb = do {
 7     use Regexp::Grammars;
 8     qr{
 9         <copy_cmd>
10
11         <rule: copy_cmd>
12             copy (<from=file>) (<to=file>) <type=(?{ 'std' })>
13             | <from=file> -> <to=file> <type=(?{ 'fwd' })>
14             | <to=file> <- <from=file> <type=(?{ 'bwd' })>
15             | .+ (?{ die "Syntax error!\n" })
16
17         <token: file> [\w./\\]+
18     }x;
19 };
20
21 while (my $input = <>) {
22     while ($input =~ m{$rbb}g) {
23         say("matches: <$&>");
24         say Dumper \% /;
25     }
26 }

```

La ejecución muestra el comportamiento del programa con tres entradas válidas y una errónea:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 aliasedcodeblock2.pl
copy a b
matches: <copy a b
>
$VAR1 = {
    '' => 'copy a b
',
    'copy_cmd' => {
        '' => 'copy a b
',
        'to' => 'b',
        'from' => 'a',
        'type' => 'std'
    }
}

```

⁸Versión de `Grammar.pm` obtenida por email con las correcciones de Damian

```

};

b <- a
matches: <b <- a
>
$VAR1 = {
    '' => 'b <- a
    ,
        'copy_cmd' => {
            '' => 'b <- a
            ,
                'to' => 'b',
                'from' => 'a',
                'type' => 'bwd'
            }
        }
};

```

```

a -> b
matches: <a -> b
>
$VAR1 = {
    '' => 'a -> b
    ,
        'copy_cmd' => {
            '' => 'a -> b
            ,
                'to' => 'b',
                'from' => 'a',
                'type' => 'fwd'
            }
        }
};

```

```

cp a b
Syntax error!

```

Pseudo subreglas y depuración

Note that, in addition to the semantics described above, aliased subpatterns and code blocks also become visible to Regexp::Grammars integrated debugger (see Debugging).

3.10.6. Llamadas a subreglas desmemoriadas

By default, every subrule call saves its result into the result-hash, either under its own name, or under an alias.

However, sometimes you may want to refactor some literal part of a rule into one or more subrules, without having those submatches added to the result-hash. The syntax for calling a subrule, but ignoring its return value is:

```
<.SUBRULE>
```

(which is stolen directly from Perl 6).

For example, you may prefer to rewrite a rule such as:

```
<rule: paren_pair>
```



```

\ (
  (?: <escape> | <paren_pair> | <brace_pair> | [^()] ) *
\ )

```

without any literal matching, like so:

```

<rule: paren_pair>

  <.left_paren>
    (?: <escape> | <paren_pair> | <brace_pair> | <.non_paren> ) *
  <.right_paren>

<token: left_paren>  \ (
<token: right_paren> \ )
<token: non_paren>   [ ^ ( ) ]

```

Moreover, as the individual components inside the parentheses probably aren't being captured for any useful purpose either, you could further optimize that to:

```

<rule: paren_pair>

  <.left_paren>
    (?: <.escape> | <.paren_pair> | <.brace_pair> | <.non_paren> ) *
  <.right_paren>

```

Note that you can also use the dot modifier on an aliased subpattern:

```

<.Alias= (SUBPATTERN) >

```

This seemingly contradictory behaviour (of giving a subpattern a name, then deliberately ignoring that name) actually does make sense in one situation. Providing the alias makes the subpattern visible to the debugger, while using the dot stops it from affecting the result-hash. See *Debugging non-grammars* for an example of this usage.

Ejemplo: Números entre comas

Por ejemplo, queremos reconocer listas de números separados por comas. Supongamos también que queremos darle un nombre a la expresión regular de separación. Quizá, aunque no es el caso, porque la expresión regular de separación sea suficientemente compleja. Si no usamos la notación *punto* la coma aparecerá en la estructura:

```

pl@nereida:~/Lregexpg grammars/demo$ cat -n numberscomma.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5 $Data::Dumper::Indent = 1;
 6
 7 my $rbb = do {
 8     use Regexp::Grammars;
 9
10     qr{
11         <numbers>
12
13         <objrule: numbers>

```

```

14         <[number]> (<comma> <[number]>)*
15
16         <objtoken: number> \s*\d+
17         <token: comma> \s*,
18     }xms;
19 };
20
21 while (my $input = <>) {
22     if ($input =~ m{$rbb}) {
23         say("matches: <${&}>");
24         say Dumper \% /;
25     }
26 }

```

En efecto, aparece la clave comma:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 numberscomma.pl
2, 3, 4
matches: <2, 3, 4>
$VAR1 = {
    '' => '2, 3, 4',
    'numbers' => bless( {
        '' => '2, 3, 4',
        'number' => [
            bless( { '' => '2' }, 'number' ),
            bless( { '' => '3' }, 'number' ),
            bless( { '' => '4' }, 'number' )
        ],
        'comma' => ', ',
    }, 'numbers' )
};

```

Si cambiamos la llamada a la regla <comma> por <.comma>

```

pl@nereida:~/Lregexpggrammars/demo$ diff numberscomma.pl numberscomma2.pl
14c14
<         <[number]> (<comma> <[number]>)*
---
>         <[number]> (<.comma> <[number]>)*

```

eliminamos la aparición de la innecesaria clave:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 numberscomma2.pl
2, 3, 4
matches: <2, 3, 4>
$VAR1 = {
    '' => '2, 3, 4',
    'numbers' => bless( {
        '' => '2, 3, 4',
        'number' => [
            bless( { '' => '2' }, 'number' ),
            bless( { '' => '3' }, 'number' ),
            bless( { '' => '4' }, 'number' )
        ],
    }, 'numbers' )
};

```

3.10.7. Destilación del resultado

Destilación manual

Regexp::Grammars also offers full manual control over the distillation process. If you use the reserved word MATCH as the alias for a subrule call:

```
<MATCH=filename>
```

or a subpattern match:

```
<MATCH=( \w+ )>
```

or a code block:

```
<MATCH=(?{ 42 })>
```

then the current rule will treat the return value of that subrule, pattern, or code block as its complete result, and return that value instead of the usual result-hash it constructs. This is the case even if the result has other entries that would normally also be returned.

For example, in a rule like:

```
<rule: term>
  <MATCH=literal>
  | <left_paren> <MATCH=expr> <right_paren>
```

The use of MATCH aliases causes the rule to return either whatever <literal> returns, or whatever <expr> returns (provided it's between left and right parentheses).

Note that, in this second case, even though <left_paren> and <right_paren> are captured to the result-hash, they are not returned, because the MATCH alias overrides the normal return the result-hash semantics and returns only what its associated subrule (i.e. <expr>) produces.

El siguiente ejemplo ilustra el uso del alias MATCH:

```
$ cat -n demo_calc.pl
1  #!/usr/local/lib/perl/5.10.1/bin/perl5.10.1
2  use v5.10;
3  use warnings;
4
5  my $calculator = do{
6      use Regexp::Grammars;
7      qr{
8          <Answer>
9
10         <rule: Answer>
11             <X=Mult> <Op=([+-])> <Y=Answer>
12             | <MATCH=Mult>
13
14         <rule: Mult>
15             <X=Pow> <Op=([*/%])> <Y=Mult>
16             | <MATCH=Pow>
17
18         <rule: Pow>
19             <X=Term> <Op=(\^)> <Y=Pow>
20             | <MATCH=Term>
```

```

21
22     <rule: Term>
23         <MATCH=Literal>
24         | \ ( <MATCH=Answer> \ )
25
26     <token: Literal>
27         <MATCH=( [+ -]? \d++ (?: \. \d++ )?+ )>
28     }xms
29 };
30
31 while (my $input = <>) {
32     if ($input =~ $calculator) {
33         use Data::Dumper 'Dumper';
34         warn Dumper \% /;
35     }
36 }

```

Veamos una ejecución:

```

$ ./demo_calc.pl
2+3*5
$VAR1 = {
    '' => '2+3*5',
    'Answer' => {
        '' => '2+3*5',
        'Op' => '+',
        'X' => '2',
        'Y' => {
            '' => '3*5',
            'Op' => '*',
            'X' => '3',
            'Y' => '5'
        }
    }
};
4-5-2
$VAR1 = {
    '' => '4-5-2',
    'Answer' => {
        '' => '4-5-2',
        'Op' => '-',
        'X' => '4',
        'Y' => {
            '' => '5-2',
            'Op' => '-',
            'X' => '5',
            'Y' => '2'
        }
    }
};

```

Obsérvese como el árbol construido para la expresión 4-5-2 se hunde a derechas dando lugar a una jerarquía errónea. Para arreglar el problema sería necesario eliminar la recursividad por la izquierda en las reglas correspondientes.

Destilación en el programa

It's also possible to control what a rule returns from within a code block. `Regexp::Grammars` provides a set of reserved variables that give direct access to the result-hash.

The result-hash itself can be accessed as `%MATCH` within any code block inside a rule. For example:

```
<rule: sum>
  <X=product> \+ <Y=product>
    <MATCH=(?{ $MATCH{X} + $MATCH{Y} })>
```

Here, the rule matches a product (aliased 'X' in the result-hash), then a literal '+', then another product (aliased to 'Y' in the result-hash). The rule then executes the code block, which accesses the two saved values (as `$MATCH{X}` and `$MATCH{Y}`), adding them together. Because the block is itself aliased to `MATCH`, the sum produced by the block becomes the (only) result of the rule.

It is also possible to set the rule result from within a code block (instead of aliasing it). The special override return value is represented by the special variable `$MATCH`. So the previous example could be rewritten:

```
<rule: sum>
  <X=product> \+ <Y=product>
    (?{ $MATCH = $MATCH{X} + $MATCH{Y} })
```

Both forms are identical in effect. Any assignment to `$MATCH` overrides the normal return all subrule results behaviour.

Assigning to `$MATCH` directly is particularly handy if the result may not always be distillable, for example:

```
<rule: sum>
  <X=product> \+ <Y=product>
    (?{ if (!ref $MATCH{X} && !ref $MATCH{Y}) {
      # Reduce to sum, if both terms are simple scalars...
      $MATCH = $MATCH{X} + $MATCH{Y};
    }
    else {
      # Return full syntax tree for non-simple case...
      $MATCH{op} = '+';
    }
  })
```

Note that you can also partially override the subrule return behaviour. Normally, the subrule returns the complete text it matched under the empty key of its result-hash. That is, of course, `$MATCH{""}`, so you can override just that behaviour by directly assigning to that entry.

For example, if you have a rule that matches key/value pairs from a configuration file, you might prefer that any trailing comments not be included in the matched text entry of the rule's result-hash. You could hide such comments like so:

```
<rule: config_line>
  <key> : <value> <comment>?
  (?{
    # Edit trailing comments out of "matched text" entry...
    $MATCH = "$MATCH{key} : $MATCH{value}";
  })
```

Some more examples of the uses of \$MATCH:

```
<rule: FuncDecl>
  # Keyword   Name                                Keep return the name (as a string)...
  func       <Identifier> ;      (?{ $MATCH = $MATCH{'Identifier'} })

<rule: NumList>
  # Numbers in square brackets...
  \[
    ( \d+ (?: , \d+)* )
  \]

  # Return only the numbers...
  (?{ $MATCH = $CAPTURE })

<token: Cmd>
  # Match standard variants then standardize the keyword...
  (?: mv | move | rename )      (?{ $MATCH = 'mv'; })
```

\$CAPTURE and \$CONTEXT are both aliases for the built-in read-only \$^N variable, which always contains the substring matched by the nearest preceding (...) capture. \$^N still works perfectly well, but these are provided to improve the readability of code blocks and error messages respectively.

El siguiente código implementa una calculadora usando destilación en el código:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n demo_calc_inline.pl
1  use v5.10;
2  use warnings;
3
4  my $calculator = do{
5      use Regexp::Grammars;
6      qr{
7          <Answer>
8
9          <rule: Answer>
10             <X=Mult> \+ <Y=Answer>
11                 (?{ $MATCH = $MATCH{X} + $MATCH{Y}; })
12             | <X=Mult> - <Y=Answer>
13                 (?{ $MATCH = $MATCH{X} - $MATCH{Y}; })
14             | <MATCH=Mult>
15
16             <rule: Mult>
17                 <X=Pow> \* <Y=Mult>
18                     (?{ $MATCH = $MATCH{X} * $MATCH{Y}; })
19                 | <X=Pow> / <Y=Mult>
20                     (?{ $MATCH = $MATCH{X} / $MATCH{Y}; })
21                 | <X=Pow> % <Y=Mult>
22                     (?{ $MATCH = $MATCH{X} % $MATCH{Y}; })
23                 | <MATCH=Pow>
24
25             <rule: Pow>
```

```

26         <X=Term> \^ <Y=Pow>
27         ({ $MATCH = $MATCH{X} ** $MATCH{Y}; })
28     | <MATCH=Term>
29
30     <rule: Term>
31         <MATCH=Literal>
32     | \ ( <MATCH=Answer> \ )
33
34     <token: Literal>
35         <MATCH=( [+ -]? \d++ (?: \. \d++ )?+ )>
36     }xms
37 };
38
39 while (my $input = <>) {
40     if ($input =~ $calculator) {
41         say '--> ', $/{Answer};
42     }
43 }

```

Ejercicio 3.10.1. *Cual es la salida del programa anterior para las entradas:*

- 4-2-2
- 8/4/2
- 2^2^3

3.10.8. Llamadas privadas a subreglas y subreglas privadas

If a rule name (or an alias) begins with an underscore:

```

    <_RULENAME>         <_ALIAS=RULENAME>
    <[_RULENAME]>       <[_ALIAS=RULENAME]>

```

then matching proceeds as normal, and any result that is returned is stored in the current result-hash in the usual way.

However, when any rule finishes (and just before it returns) it first filters its result-hash, removing any entries whose keys begin with an underscore. This means that any subrule with an underscored name (or with an underscored alias) remembers its result, but only until the end of the current rule. Its results are effectively private to the current rule.

This is especially useful in conjunction with result distillation.

3.10.9. Mas sobre listas

Reconocimiento manual de listas

Analizando listas manualmente

El siguiente ejemplo muestra como construir un reconocedor de listas (posiblemente vacías) de números:

```

casiano@millo:~/Lregex-grammar-examples$ cat -n simple_list.pl
1  #!/soft/perl5lib/bin/perl5.10.1
2  use v5.10;
3
4  use Regexp::Grammars;
5

```

```

6 my $list = qr{
7     <List>
8
9     <rule: List>
10         <digit> <List>
11         | # empty
12
13     <rule: digit>
14         <MATCH=(\d+)>
15
16 }xms;
17
18 while (my $input = <>) {
19     chomp $input;
20     if ($input =~ $list) {
21         use Data::Dumper 'Dumper';
22         warn Dumper \% /;
23     }
24     else {
25         warn "Does not match\n"
26     }
27 }

```

Sigue una ejecución:

```
casiano@millo:~/Lregex-grammar-examples$ ./simple_list.pl
```

```
2 3 4
```

```

$VAR1 = {
    '' => '2 3 4',
    'List' => {
        '' => '2 3 4',
        'digit' => '2'
        'List' => {
            '' => '3 4',
            'digit' => '3'
            'List' => {
                '' => '4',
                'digit' => '4'
                'List' => '',
            },
        },
    },
};

```

Influencia del orden en el lenguaje reconocido

Tenga en cuenta que el orden de las reglas influye en el lenguaje reconocido. Véase lo que ocurre si cambiamos en el ejemplo anterior el orden de las reglas:

```
casiano@millo:~/Lregex-grammar-examples$ cat -n simple_list_empty_first.pl
```

```

1  #!/soft/perl5lib/bin/perl5.10.1
2  use v5.10;
3
4  use Regexp::Grammars;
5
6  my $list = qr{

```



```

7      <List>
8
9      <rule: List>
10         # empty
11         | <digit> <List>
12
13     <rule: digit>
14         <MATCH=(\d+)>
15
16 }xms;
17
18 while (my $input = <>) {
19     chomp $input;
20     if ($input =~ $list) {
21         use Data::Dumper 'Dumper';
22         warn Dumper \% /;
23     }
24     else {
25         warn "Does not match\n"
26     }
27 }

```

Al ejecutar se obtiene:

```

casiano@millo:~/Lregex-grammar-examples$ ./simple_list_empty_first.pl
2 3 4
$VAR1 = {
    '' => '',
    'List' => ''
};

```

Por supuesto basta poner anclas en el patrón a buscar para forzar a que se reconozca la lista completa:

```

pl@nereida:~/Lregexgrammars/demo$ diff simple_list_empty_first.pl simple_list_empty_first_with_anchors.pl
7c7
<      <List>
---
>      ^<List>$

```

En efecto, la nueva versión reconoce la lista:

```

pl@nereida:~/Lregexgrammars/demo$ perl5.10.1 simple_list_empty_first_with_anchors.pl
2 3 4
$VAR1 = {
    '' => '2 3 4',
    'List' => {
        'List' => {
            'List' => '',
            '' => '4',
            'digit' => '4'
        },
        '' => '3 4',
        'digit' => '3'
    },
};

```

```

        '' => '2 3 4',
        'digit' => '2'
    }
};

```

Si se quiere mantener la producción vacía en primer lugar pero forzar el reconocimiento de la lista completa, se puede hacer uso de un lookahead negativo:

```

pl@nereida:~/Lregexpggrammars/demo$ cat -n simple_list_empty_first_with_lookahead.pl
 1  #!/soft/perl5lib/bin/perl5.10.1
 2  use v5.10;
 3
 4  use strict;
 5  use Regexp::Grammars;
 6
 7  my $list = qr{
 8      <List>
 9
10      <rule: List>
11          (?! <digit> ) # still empty production
12          | <digit> <List>
13
14      <rule: digit>
15          <MATCH=(\d+)>
16
17  }xms;
18
19  while (my $input = <>) {
20      chomp $input;
21      if ($input =~ $list) {
22          use Data::Dumper 'Dumper';
23          warn Dumper \% /;
24      }
25      else {
26          warn "Does not match\n"
27      }
28  }

```

Así, sólo se reducirá por la regla vacía si el siguiente token no es un número. Sigue un ejemplo de ejecución:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 simple_list_empty_first_with_lookahead.pl
2 3 4
$VAR1 = {
    '' => '2 3 4',
    'List' => {
        'List' => {
            'List' => '',
            '' => '4',
            'digit' => '4',
        },
        '' => '3 4',
        'digit' => '3',
    },
};

```

```

        '' => '2 3 4',
        'digit' => '2'
    }
};

```

Aplanamiento manual de listas

¿Cómo podemos hacer que la estructura retornada por el reconocedor sea una lista?. Podemos añadir acciones como sigue:

```
casiano@millor:~/Lregex-grammar-examples$ cat -n simple_list_action.pl
```

```

1  #!/soft/perl5lib/bin/perl5.10.1
2  use v5.10;
3
4  use Regexp::Grammars;
5
6  my $list = qr{
7      <List>
8
9      <rule: List>
10         <digit> <X=List> <MATCH= (?{ unshift @{$MATCH{X}}, $MATCH{digit}; $MATCH{X} }
11         | # empty
12         <MATCH= (?{ [] })>
13
14     <rule: digit>
15         <MATCH=(\d+)>
16
17 }xms;
18
19 while (my $input = <>) {
20     chomp $input;
21     if ($input =~ $list) {
22         use Data::Dumper 'Dumper';
23         warn Dumper \% /;
24     }
25     else {
26         warn "Does not match\n"
27     }
28 }

```

Al ejecutarse este programa produce una salida como:

```

pl@nereida:~/Lregexgrammars/demo$ perl5.10.1 simple_list_action.pl
2 3 4
$VAR1 = {
    '' => '2 3 4',
    'List' => [ '2', '3', '4' ]
};

```

Los operadores de repetición

Los operadores de repetición como *, +, etc. permiten simplificar el análisis de lenguajes de listas:

```

pl@nereida:~/Lregexgrammars/demo$ cat -n simple_list_star.pl
1  #!/soft/perl5lib/bin/perl5.10.1

```

```

2  use v5.10;
3
4  use Regexp::Grammars;
5
6  my $list = qr{
7      <List>
8
9      <rule: List>
10         (?: <[digit]>)*
11
12         <rule: digit>
13             <MATCH=(\d+)>
14
15 }xms;
16
17 while (my $input = <>) {
18     chomp $input;
19     if ($input =~ $list) {
20         use Data::Dumper 'Dumper';
21         warn Dumper \% /;
22     }
23     else {
24         warn "Does not match\n"
25     }
26 }

```

Los corchetes alrededor de `digit` hacen que el valor asociado con el patrón sea la lista de números. Si no los ponemos el valor asociado sería el último valor de la lista.

Listas separadas por Algo

One of the commonest tasks in text parsing is to match a list of unspecified length, in which items are separated by a fixed token. Things like:

```

1, 2, 3 , 4 ,13, 91      # Numbers separated by commas and spaces

g-c-a-g-t-t-a-c-a      # Bases separated by dashes

/usr/local/bin          # Names separated by directory markers

/usr:/usr/local:bin     # Directories separated by colons

```

The usual construct required to parse these kinds of structures is either:

```

<rule: list>

    <item> <separator> <list>      # recursive definition
| <item>                          # base case

```

Or, more efficiently, but less prettily:

```

<rule: list>

    <[item]> (?: <separator> <[item]> )*  # iterative definition

```

Because this is such a common requirement, `Regexp::Grammars` provides a cleaner way to specify the iterative version. The syntax is taken from Perl 6:

```
<rule: list>
```

```
    <[item]> ** <separator>                                # iterative definition
```

*This is a repetition specifier on the first subrule (hence the use of `**` as the marker, to reflect the repetitive behaviour of `*`). However, the number of repetitions is controlled by the second subrule: the first subrule will be repeatedly matched for as long as the second subrule matches immediately after it.*

So, for example, you can match a sequence of numbers separated by commas with:

```
<[number]> ** <comma>
```

```
<token: number>  \d+
<token: comma>   \s* , \s*
```

Note that it's important to use the `<[...]>` form for the items being matched, so that all of them are saved in the result hash. You can also save all the separators (if that's important):

```
<[number]> ** <[comma]>
```

The repeated item must be specified as a subrule call for some kind, but the separators may be specified either as a subrule or a bracketed pattern. For example:

```
<[number]> ** ( , )
```

The separator must always be specified in matched delimiters of some kind: either matching `<...>` or matching `(...)`. A common error is to write:

```
<[number]> ** ,
```

You can also use a pattern as the item matcher, but it must be aliased into a subrule:

```
<[item=(\d+)]> ** ( , )
```

Ejemplo: Listas de números separados por comas

Veamos un ejemplo sencillo:

```
casiano@millo:~/src/perl/regexp-grammar-examples$ cat -n demo_list.pl
 1  #!/soft/perl5lib/bin/perl5.10.1
 2  use v5.10;
 3
 4  use Regexp::Grammars;
 5
 6  my $list_nonempty = qr{
 7      <List>
 8
 9      <rule: List>
10          \( <[Value]> ** (,) \)
11
12      <token: Value>
13          \d+
```

```

14 }xms;
15
16 my $list_empty = qr{
17     <List>
18
19     <rule: List>
20         \(\ (?<[Value]> ** <_Sep=(,)> )? \)
21
22     <token: Value>
23         \d+
24 }xms;
25
26 use Smart::Comments;
27
28
29 while (my $input = <>) {
30     my $input2 = $input;
31     if ($input =~ $list_nonempty) {
32         ### nonempty: $/{List}
33     }
34     if ($input2 =~ $list_empty) {
35         ### empty: $/{List}
36     }
37 }

```

Sigue un ejemplo de ejecución:

```

casiano@millo:~/src/perl/regexp-grammar-examples$ ./demo_list.pl
(3,4,5)

```

```

### nonempty: {
###     '' => '(3,4,5)',
###     Value => [
###         '3',
###         '4',
###         '5'
###     ]
### }

```

```

### empty: {
###     '' => '(3,4,5)',
###     Value => [
###         '3',
###         '4',
###         '5'
###     ]
### }
(

```

```

### empty: '()'

```

Ejemplo: AST para las expresiones aritméticas

Las expresiones aritméticas puede definirse como una jerarquía de listas como sigue:

```

pl@nereida:~/Lregexpggrammars/demo$ cat -n calcaslist.pl
 1 use strict;
 2 use warnings;
 3 use 5.010;
 4 use Data::Dumper;
 5 $Data::Dumper::Indent = 1;
 6
 7 my $rbb = do {
 8     use Regexp::Grammars;
 9
10     qr{
11         \A<expr>\z
12
13         <objrule: expr>      <[operands=term]> ** <[operators=addop]>
14
15         <objrule: term>      <[operands=uneg]> ** <[operators=mulop]>
16
17         <objrule: uneg>      <[operators=minus]>* <[operands=power]>
18
19         <objrule: power>     <[operands=factorial]> ** <[operators=powerop]>
20
21         <objrule: factorial> <[operands=factor]> <[operators=(!)]>*
22
23         <objrule: factor>    <val=( [+ - ] ? \d + ( ? : \. \d * ) ? ) >
24                             | \ ( <MATCH=expr> \ )
25
26         <token: addop>       [ + - ]
27
28         <token: mulop>       [ * / ]
29
30         <token: powerop>     \ * \ * | \ ^
31
32         <token: minus>       - <MATCH=( ? { 'NEG' } ) >
33
34     }x;
35 };
36
37 while (my $input = <>) {
38     chomp($input);
39     if ($input =~ m{$rbb}) {
40         my $tree = $/{expr};
41         say Dumper $tree;
42     }
43     else {
44         say("does not match");
45     }
46 }
47 }

```

Obsérvese el árbol generado para la expresión 4-2-2:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 calcaslist.pl
4-2-2
$VAR1 = bless( {

```

```

operands' => [
  bless( {
    'operands' => [
      bless( {
        'operands' => [
          bless( {
            'operands' => [
              bless( { '' => '4', 'val' => '4' }, 'factor' )
            ],
            '' => '4',
          }, 'factorial' )
        ],
        '' => '4',
      }, 'power' )
    ],
    '' => '4',
  }, 'uneg' )
],
'' => '4',
}, 'term' ),
bless( {
  'operands' => [
    bless( {
      'operands' => [
        bless( {
          'operands' => [
            bless( { '' => '2', 'val' => '2' }, 'factor' )
          ],
          '' => '2',
        }, 'factorial' )
      ],
      '' => '2',
    }, 'power' )
  ],
  '' => '2',
}, 'uneg' )
],
'' => '2',
}, 'term' ),
bless( {
  'operands' => [
    bless( {
      'operands' => [
        bless( {
          'operands' => [
            bless( { '' => '2', 'val' => '2' }, 'factor' )
          ],

```



```

        '' => '2'
    }, 'factorial' )
],
    '' => '2'
}, 'power' )
],
    '' => '2'
}, 'uneg' )
],
    '' => '2'
}, 'term' )
],
'' => '4-2-2',
'operators' => [
    '-',
    '-'
]
}, 'expr' );

```

3.10.10. La directiva require

La directiva `require` es similar en su funcionamiento al paréntesis 5.10 (`??{ Código Perl }`) el cuál hace que el Código Perl sea evaluado durante el tiempo de matching. El resultado de la evaluación se trata como una expresión regular con la que deberá casarse. (véase la sección 3.2.9 para mas detalles).

La sintáxis de la directiva `<require:>` es

```
<require: (?{ CODE }) >
```

The code block is executed and if its final value is true, matching continues from the same position. If the block's final value is false, the match fails at that point and starts backtracking.

The `<require:...>` directive is useful for testing conditions that it's not easy (or even possible) to check within the syntax of the the regex itself. For example:

```

<rule: IPV4_Octet_Decimal>
    # Up three digits...
    <MATCH= ( \d{1,3}+ )>

    # ...but less than 256...
    <require: (?{ $MATCH <= 255 })>

```

A require expects a regex codeblock as its argument and succeeds if the final value of that codeblock is true. If the final value is false, the directive fails and the rule starts backtracking.

Note, in this example that the digits are matched with `\d{1,3}+`. The trailing `+` prevents the `{1,3}` repetition from backtracking to a smaller number of digits if the `<require:...>` fails.

El programa `demo_IP4.pl` ilustra el uso de la directiva:

```

pl@nereida:~/Lregexpgrammars/demo$ cat -n ./demo_IP4.pl
1  #!/usr/bin/env perl5.10.1
2  use v5.10;
3  use warnings;

```

```

4
5 use Regexp::Grammars;
6
7 my $grammar = qr{
8     \A <IP4_addr> \Z
9
10    <token: quad>
11        <MATCH=(\d{1,3})>
12        <require: (?{ $MATCH < 256 })>
13
14    <token: IP4_addr>
15        <[MATCH=quad]>*(\.)
16        <require: (?{ @$MATCH == 4 })>
17 }xms;
18
19 while (my $line = <>) {
20     if ($line =~ $grammar) {
21         use Data::Dumper 'Dumper';
22         say Dumper \%/;
23     }
24     else {
25         say 'Does not match'
26     }
27 }

```

Las condiciones usadas en el `require` obligan a que cada quad⁹ sea menor que 256 y a que existan sólo cuatro quads.

Sigue un ejemplo de ejecución:

```

pl@nereida:~/Lregexpggrammars/demo$ ./demo_IP4.pl
123 . 145 . 105 . 252
Does not match
pl@nereida:~/Lregexpggrammars/demo$ ./demo_IP4.pl
123.145.105.252
$VAR1 = {
    '' => '123.145.105.252',
    'IP4_addr' => [
        123,
        145,
        105,
        252
    ]
};
pl@nereida:~/Lregexpggrammars/demo$ ./demo_IP4.pl
148.257.128.128
Does not match
0.0.0.299
Does not match
pl@nereida:~/Lregexpggrammars/demo$ ./demo_IP4.pl
123.145.105.242.193

```

⁹ A quad (pronounced KWAHD) is a unit in a set of something that comes in four units. The term is sometimes used to describe each of the four numbers that constitute an Internet Protocol (IP) address. Thus, an Internet address in its numeric form (which is also sometimes called a dot address) consists of four quads separated by "dots"(periods).

A quad also means *a quarter* in some usages. (A quarter as a U.S. coin or monetary unit means *a quarter of a dollar*, and in slang is sometimes called *two bits*. However, this usage does not mean two binary bits as used in computers.)

Does not match

Obsérvese como no se aceptan blancos entre los puntos en esta versión. ¿Sabría explicar la causa?

3.10.11. Casando con las claves de un hash

In some situations a grammar may need a rule that matches dozens, hundreds, or even thousands of one-word alternatives. For example, when matching command names, or valid userids, or English words. In such cases it is often impractical (and always inefficient) to list all the alternatives between | alterators:

```
<rule: shell_cmd>
  a2p | ac | apply | ar | automake | awk | ...
  # ...and 400 lines later
  ... | zdiff | zgrep | zip | zmore | zsh

<rule: valid_word>
  a | aa | aal | aalii | aam | aardvark | aardwolf | aba | ...
  # ...and 40,000 lines later...
  ... | zymotize | zymotoxic | zymurgy | zythem | zythum
```

To simplify such cases, Regexp::Grammars provides a special construct that allows you to specify all the alternatives as the keys of a normal hash. The syntax for that construct is simply to put the hash name inside angle brackets (with no space between the angles and the hash name).

Which means that the rules in the previous example could also be written:

```
<rule: shell_cmd>
  <%cmds>

<rule: valid_word>
  <%dict>
```

provided that the two hashes (%cmds and %dict) are visible in the scope where the grammar is created.

Internally, the construct is converted to something equivalent to:

```
<rule: shell_cmd>
  (<.hk>) <require: exists $cmds{$CAPTURE}>

<rule: valid_word>
  (<.hk>) <require: exists $dict{$CAPTURE}>
```

The special <hk> rule is created automatically, and defaults to \S+, but you can also define it explicitly to handle other kinds of keys. For example:

```
<rule: hk>
  .+          # Key may be any number of chars on a single line

<rule: hk>
  [ACGT]{10,} # Key is a base sequence of at least 10 pairs
```

Matching a hash key in this way is typically significantly faster than matching a full set of alternations. Specifically, it is $O(\text{length of longest potential key})$, instead of $O(\text{number of keys})$.

Ejemplo de uso de la directiva hash

Sigue un ejemplo:

```
pl@nereida:~/Lregexpggrammars/demo$ cat -n hash.pl
1  #!/usr/bin/env perl5.10.1
2  use strict;
3  use warnings;
4  use 5.010;
5  use Data::Dumper;
6  $Data::Dumper::Deparse = 1;
7
8  my %cmd = map { ($_ => undef ) } qw( uname pwd date );
9
10 my $rbb = do {
11     use Regexp::Grammars;
12
13     qr{
14         ^<command>$
15
16         <rule: command>
17             <cmd=%cmd> (?: <[arg]> )*
18
19             <token: arg> [^\s<>'&]+
20     }xms;
21 };
22
23 while (my $input = <>) {
24     chomp($input);
25     if ($input =~ m{$rbb}) {
26         say("matches: <$&>");
27         say Dumper \% /;
28         system $/{''}
29     }
30     else {
31         say("does not match");
32     }
33 }
```

Sigue un ejemplo de ejecución:

```
pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 hash.pl
a2p f1 f2
matches: <a2p f1 f2>
$VAR1 = {
    '' => 'a2p f1 f2',
    'command' => {
        '' => 'a2p f1 f2',
        'cmd' => 'a2p',
        'arg' => [
            'f1',
            'f2'
        ]
    }
};
```

```
pocho 2 5
does not match
```

3.10.12. Depuración

Regexp::Grammars provides a number of features specifically designed to help debug both grammars and the data they parse.

All debugging messages are written to a log file (which, by default, is just STDERR). However, you can specify a disk file explicitly by placing a "<logfile:..." directive at the start of your grammar¹⁰:

```
$grammar = qr{

    <logfile: LaTeX_parser_log >

    \A <LaTeX_file> \Z      # Pattern to match

    <rule: LaTeX_file>
        # etc.
}x;
```

You can also explicitly specify that messages go to the terminal:

```
<logfile: - >
```

Debugging grammar creation

Whenever a log file has been directly specified, Regexp::Grammars automatically does verbose static analysis of your grammar. That is, whenever it compiles a grammar containing an explicit "<logfile:..." directive it logs a series of messages explaining how it has interpreted the various components of that grammar. For example, the following grammar:

```
pl@nereida:~/Lregexgrammars/demo$ cat -n log.pl
```

```
1  #!/usr/bin/env perl5.10.1
2  use strict;
3  use warnings;
4  use 5.010;
5  use Data::Dumper;
6
7  my $rbb = do {
8      use Regexp::Grammars;
9
10     qr{
11         <logfile: ->
12
13         <numbers>
14
15         <rule: numbers>
16             <number> ** <comma>
17
18         <token: number> \d+
19     }
```

¹⁰no funcionará si no se pone al principio de la gramática

```

20         <token: comma>    ,
21     }xms;
22 };
23
24 while (my $input = <>) {
25     if ($input =~ m/${rbb}) {
26         say("matches: <$&>");
27         say Dumper \% /;
28     }
29 }

```

would produce the following analysis in the terminal:

```

pl@nereida:~/Lregexgrammars/demo$ ./log.pl
warn | Repeated subrule <number>* will only capture its final match
      | (Did you mean <[number]>* instead?)
      |
info  | Processing the main regex before any rule definitions
      |
      | ...Treating <numbers> as:
      | | match the subrule <numbers>
      | | \ saving the match in $MATCH{'numbers'}
      | |
      | | \___End of main regex
      |
      | Defining a rule: <numbers>
      | | ...Returns: a hash
      | |
      | | ...Treating <number> as:
      | | | match the subrule <number>
      | | | \ saving the match in $MATCH{'number'}
      | |
      | | ...Treating <.comma> as:
      | | | match the subrule <comma>
      | | | \ but don't save anything
      | |
      | | ...Treating <number> ** <.comma> as:
      | | | repeatedly match the subrule <number>
      | | | \ as long as the matches are separated by matches of <.comma>
      | |
      | | \___End of rule definition
      |
      | Defining a rule: <number>
      | | ...Returns: a hash
      | |
      | | ...Treating '\d' as:
      | | | \ normal Perl regex syntax
      | |
      | | ...Treating '+ ' as:
      | | | \ normal Perl regex syntax
      | |
      | | \___End of rule definition
      |
      | Defining a rule: <comma>

```

```

|      |...Returns: a hash
|      |
|      |...Treating ', ' as:
|      |      \ normal Perl regex syntax
|      |
|      |___End of rule definition
|
2, 3, 4
matches: <2, 3, 4>
$VAR1 = {
    '' => '2, 3, 4',
    'numbers' => {
        '' => '2, 3, 4',
        'number' => '4'
    }
};

```

This kind of static analysis is a useful starting point in debugging a miscreant grammar¹¹, because it enables you to see what you actually specified (as opposed to what you thought you'd specified).

Debugging grammar execution

Regexp::Grammars also provides a simple interactive debugger, with which you can observe the process of parsing and the data being collected in any result-hash.

To initiate debugging, place a <debug:...> directive anywhere in your grammar. When parsing reaches that directive the debugger will be activated, and the command specified in the directive immediately executed. The available commands are:

```

<debug: on>      - Enable debugging, stop when entire grammar matches
<debug: match>   - Enable debugging, stop when a rule matches
<debug: try>     - Enable debugging, stop when a rule is tried
<debug: off>     - Disable debugging and continue parsing silently

<debug: continue> - Synonym for <debug: on>
<debug: run>      - Synonym for <debug: on>
<debug: step>     - Synonym for <debug: try>

```

These directives can be placed anywhere within a grammar and take effect when that point is reached in the parsing. Hence, adding a <debug:step> directive is very much like setting a breakpoint at that point in the grammar. Indeed, a common debugging strategy is to turn debugging on and off only around a suspect part of the grammar:

```

<rule: tricky>   # This is where we think the problem is...
    <debug:step>
    <preamble> <text> <postscript>
    <debug:off>

```

Once the debugger is active, it steps through the parse, reporting rules that are tried, matches and failures, backtracking and restarts, and the parser's location within both the grammar and the text being matched. That report looks like this:

¹¹ miscreant - One who has behaved badly, or illegally; One not restrained by moral principles; an unscrupulous villain; One who holds an incorrect religious belief; an unbeliever; Lacking in conscience or moral principles; unscrupulous; Holding an incorrect religious belief.

```

=====> Trying <grammar> from position 0
> cp file1 file2 |...Trying <cmd>
|   |...Trying <cmd=(cp)>
|   |   \FAIL <cmd=(cp)>
|   |   \FAIL <cmd>
|   |   \FAIL <grammar>
=====> Trying <grammar> from position 1
cp file1 file2 |...Trying <cmd>
|   |...Trying <cmd=(cp)>
file1 file2   |   |   \_____<cmd=(cp)> matched 'cp'
file1 file2   |   |...Trying <[file]>+
file2         |   |   \_____<[file]>+ matched 'file1'
|   |...Trying <[file]>+
[eos]         |   |   \_____<[file]>+ matched ' file2'
|   |...Trying <[file]>+
|   |   \FAIL <[file]>+
|   |...Trying <target>
|   |   |...Trying <file>
|   |   |   \FAIL <file>
|   |   |   \FAIL <target>
<~~~~~>      |   |...Backtracking 5 chars and trying new match
file2         |   |...Trying <target>
|   |   |...Trying <file>
|   |   |   \_____ <file> matched 'file2'
[eos]         |   |   \_____<target> matched 'file2'
|   |   |   \_____<cmd> matched ' cp file1 file2'
|   |   |   \_____<grammar> matched ' cp file1 file2'

```

The first column indicates the point in the input at which the parser is trying to match, as well as any backtracking or forward searching it may need to do. The remainder of the columns track the parser's hierarchical traversal of the grammar, indicating which rules are tried, which succeed, and what they match.

Provided the logfile is a terminal (as it is by default), the debugger also pauses at various points in the parsing process—before trying a rule, after a rule succeeds, or at the end of the parse—according to the most recent command issued. When it pauses, you can issue a new command by entering a single letter:

```

m      - to continue until the next subrule matches
t or s - to continue until the next subrule is tried
r or c - to continue to the end of the grammar
o      - to switch off debugging

```

Note that these are the first letters of the corresponding <debug:...> commands, listed earlier. Just hitting *ENTER* while the debugger is paused repeats the previous command.

While the debugger is paused you can also type a *d*, which will display the result-hash for the current rule. This can be useful for detecting which rule isn't returning the data you expected.

Veamos un ejemplo. El siguiente programa activa el depurador:

```

pl@nereida:~/Lregexpggrammars/demo$ cat -n demo_debug.pl
1  #!/usr/bin/env perl5.10.1
2  use 5.010;
3  use warnings;
4

```



```

5      use Regexp::Grammars;
6
7      my $balanced_brackets = qr{
8          <debug:on>
9
10         <left_delim=( \ ( )>
11         (?:
12             <[escape=( \\ )]>
13             | <recurse=( (?R) )>
14             | <[simple=( . )]>
15         )*
16         <right_delim=( \ )>
17     }xms;
18
19     while (<>) {
20         if (/$balanced_brackets/) {
21             say 'matched: ';
22             use Data::Dumper 'Dumper';
23             warn Dumper \% /;
24         }
25     }

```

Al ejecutar obtenemos

```

pl@nereida:~/Lregexpggrammars/demo$ ./demo_debug.pl
(a)
=====> Trying <grammar> from position 0
(a)\n |...Trying <left_delim=( \ ( )>
a)\n | \_____<left_delim=( \ ( )> matched '('          c
|...Trying <[escape=( \ )]>
| \FAIL <[escape=( \ )]>
|...Trying <recurse=( (?R) )>
=====> Trying <grammar> from position 1
a)\n | |...Trying <left_delim=( \ ( )>
| | \FAIL <left_delim=( \ ( )>
\FAIL <grammar>
|...Trying <[simple=( . )]>
)\n | \_____<[simple=( . )]> matched 'a'
|...Trying <[escape=( \ )]>
| \FAIL <[escape=( \ )]>
|...Trying <recurse=( (?R) )>
=====> Trying <grammar> from position 2
)\n | |...Trying <left_delim=( \ ( )>
| | \FAIL <left_delim=( \ ( )>
\FAIL <grammar>
|...Trying <[simple=( . )]>
\n | \_____<[simple=( . )]> matched ')'
|...Trying <[escape=( \ )]>
| \FAIL <[escape=( \ )]>
|...Trying <recurse=( (?R) )>
=====> Trying <grammar> from position 3
\n | |...Trying <left_delim=( \ ( )>
| | \FAIL <left_delim=( \ ( )>
\FAIL <grammar>

```

```

|...Trying <[simple=( . )]>
[eos] | \_____<[simple=( . )]> matched ''
|...Trying <[escape=( \ )]>
| \FAIL <[escape=( \ )]>
|...Trying <recurse=( (?R) )>
====> Trying <grammar> from position 4
[eos] | |...Trying <left_delim=( \ ( )>
| | \FAIL <left_delim=( \ ( )>
| \FAIL <grammar>
|...Trying <[simple=( . )]>
| \FAIL <[simple=( . )]>
|...Trying <right_delim=( \ )>
| \FAIL <right_delim=( \ )>
<~~~~~|...Backtracking 1 char and trying new match
\n |...Trying <right_delim=( \ )>
| \FAIL <right_delim=( \ )>
<~~~~~|...Backtracking 1 char and trying new match
)\n |...Trying <right_delim=( \ )>
\n | \_____<right_delim=( \ )> matched ''
| \_____<grammar> matched '(a)' d
: {
: '' => '(a)',
: 'left_delim' => '(',
: 'simple' => [
: 'a'
: ],
: 'right_delim' => ')',
: }; o
matched:
$VAR1 = {
'' => '(a)',
'left_delim' => '(',
'simple' => [
'a'
],
'right_delim' => ')',
};

```

3.10.13. Mensajes de log del usuario

Both static and interactive debugging send a series of predefined log messages to whatever log file you have specified. It is also possible to send additional, user-defined messages to the log, using the "<log:...>" directive.

This directive expects either a simple text or a codeblock as its single argument. If the argument is a code block, that code is expected to return the text of the message; if the argument is anything else, that something else is the literal message. For example:

```

<rule: ListElem>

    <Elem= ( [a-z]\d+ ) >
        <log: Checking for a suffix, too...>

    <Suffix= ( : \d+ ) >?
        <log: (?{ "ListElem: $MATCH{Elem} and $MATCH{Suffix}" })>

```

User-defined log messages implemented using a codeblock can also specify a severity level. If the codeblock of a `<log:...>` directive returns two or more values, the first is treated as a log message severity indicator, and the remaining values as separate lines of text to be logged. For example:

```
<rule: ListElem>
  <Elem= ( [a-z]\d+) >
  <Suffix= ( : \d+ ) >?

  <log: (?{
    warn => "Elem was: $MATCH{Elem}",
           "Suffix was $MATCH{Suffix}",
  })>
```

When they are encountered, user-defined log messages are interspersed between any automatic log messages (i.e. from the debugger), at the correct level of nesting for the current rule.

3.10.14. Depuración de Regexp

It is possible to use *Regexp::Grammars* without creating any subrule definitions, simply to debug a recalcitrant regex. For example, if the following regex wasn't working as expected:

```
my $balanced_brackets = qr{
  \(          # left delim
  (?:
    \\        # escape or
    | ( ?R)   # recurse or
    | .       # whatever
  )*
  \)          # right delim
}xms;
```

you could instrument it with aliased subpatterns and then debug it step-by-step, using *Regexp::Grammars*:

```
use Regexp::Grammars;

my $balanced_brackets = qr{
  <debug:step>

  <.left_delim= ( \( )>
  (?:
    <.escape= ( \\ )>
    | <.recurse= ( ?R )>
    | <.whatever=( . )>
  )*
  <.right_delim= ( \) )>
}xms;

while (<>) {
  say 'matched' if /$balanced_brackets/;
}
```

Note the use of amnesiac aliased subpatterns to avoid needlessly building a result-hash. Alternatively, you could use listifying aliases to preserve the matching structure as an additional debugging aid:

```
use Regexp::Grammars;

my $balanced_brackets = qr{
    <debug:step>

    <[left_delim= ( \ ( ) ]>
    (?
        <[escape= ( \\ ) ]>
    |   <[recurse= ( (?R) ) ]>
    |   <[whatever=( . ) ]>
    )*
    <[right_delim= ( \ ) ]>
}xms;

if ( '(a(bc)d)' =~ /$balanced_brackets/) {
    use Data::Dumper 'Dumper';
    warn Dumper \% /;
}
```

3.10.15. Manejo y recuperación de errores

En este punto debo decir que no he podido reproducir el comportamiento de las directivas `<error:>` y `<warning:>` tal y como las describe Conway en el manual de `Regexp::Grammars`.

El siguiente ejemplo ilustra un conjunto de técnicas de gestión de errores que son independientes del soporte dado por `Regexp::Grammars`.

Se trata de la misma calculadora explicada en la sección 3.10.18.

```
pl@nereida:~/Lregexpggrammars/demo/calculator$ cat -n calculatorwitherrmanagement.pl
 1  #!/usr/bin/env perl5.10.1
 2  use strict;
 3  use warnings;
 4  use 5.010;
 5  use Lingua::EN::Inflect qw(PL);
 6  use Scalar::Util qw{blessed};
 7
 8  my $rbb = do {
 9      my ($warnings, $errors);    # closure
10      sub warnings { $warnings }  # accessor
11      sub errors { $errors }      # accessor
12
13      use Regexp::Grammars;
14      qr{
15          (?{
16              $warnings = 0;
17              $errors = 0;
18          })
19          \A<expr>
20          (?
21              \z
22              |
23              (.*?) (?{
```

```

23             # Accept the string but emit a warning
24             $warnings++;
25             local our $expr = \${MATCH{expr}}{''};
26             local our $endlegal = length($$expr) > 4? "... ".substr($$expr, -4)
27             warn "Warning: Unexpected '". substr($^N, 0, 10)."'" after '$endlegal
28         })
29     )
30
31     <objrule: expr>         <[operands=term]> ** <[operators=addop]>
32
33     <objrule: term>        <[operands=uneg]> ** <[operators=mulop]>
34
35     <objrule: uneg>        <[operators=minus]>* <[operands=power]>
36
37     <objrule: power>       <[operands=factorial]> ** <[operators=powerop]>
38
39     <objrule: factorial> <[operands=factor]> <[operators=(!)]>*
40
41     <objrule: factor>      (<val=([+-]?\d+(?:\.\d*)?)>)
42                          | \(<MATCH=expr> \)
43                          | ([^~+(0-9)]+) (?{
44                              # is + and not * to avoid infinite recursion
45                              warn "Error: expecting a number or a open parent
46                              $warnings++;
47                              $errors++;
48                          }) <MATCH=factor>
49
50     <token: addop>         [+~]
51
52     <token: mulop>         [*/]
53
54     <token: powerop>       \*\*|\^
55
56     <token: minus>         - <MATCH=(?{ 'NEG' })>
57
58     }x;
59 };
60
61 sub test_calc {
62     my $prompt = shift;
63
64     print $prompt;
65     while (my $input = <>) {
66         chomp($input);
67
68         local %/;
69         $input =~ m/${rbb};
70
71         say warnings." ".PL('warning',warnings) if warnings;
72         say errors." ".PL('error',errors)         if errors;
73
74         my $tree = $/{expr};
75         if (blessed($tree)) {

```

```

76         do "PostfixCalc.pm";
77         say "postfix: ".$tree->ceval;
78
79         do "EvalCalc.pm";
80         say "result: ".$tree->ceval;
81     }
82     print $prompt;
83 }
84 say "Bye!"
85 }
86
87 ##### main
88 test_calc(
89     'Parsing infix arithmetic expressions (CTRL-D to end in unix) ',
90 );

```

Veamos algunas ejecuciones que incluyen entradas erróneas:

```

pl@nereida:~/Lregexpgrammars/demo/calculator$ ./calculatorwitherrmanagement.pl
Parsing infix arithmetic expressions (CTRL-D to end in unix) 2+3
postfix: 2 3 +
result: 5
Parsing infix arithmetic expressions (CTRL-D to end in unix) 2*(3+#)
Error: expecting a number or a open parenthesis, found: '#'
Error: expecting a number or a open parenthesis, found: '#'
Error: expecting a number or a open parenthesis, found: ')'
Warning: Unexpected '*'(3+#)' after '2'
4 warnings
3 errors
postfix: 2
result: 2
Parsing infix arithmetic expressions (CTRL-D to end in unix) 2+##4
Error: expecting a number or a open parenthesis, found: '##'
1 warning
1 error
postfix: 2 4 +
result: 6
Parsing infix arithmetic expressions (CTRL-D to end in unix) Bye!

```

Obsérvese los mensajes de error repetidos para la entrada `2*(3+#)`. Ellos son debidos a los reiterados intentos de casar `<factor>` en la regla de recuperación de errores:

```

41     <objrule: factor>    (<val=([+-]?\d+(?:\.\d*)?)>)
42                         | \(<MATCH=expr> \)
43                         | ([^-(0-9)]+) (?{
44                             # is + and not * to avoid infinite recursion
45                             warn "Error: expecting a number or a open parent
46                             $warnings++;
47                             $errors++;
48                         }) <MATCH=factor>

```

en este caso resulta imposible encontrar un factor. Se puede cambiar la conducta indicando un `(* COMMIT)` antes de la llamada a `<MATCH=factor>`:

```

41     <objrule: factor>    (<val=([+-]?\d+(?:\.\d*)?)>)
42                         | \(<MATCH=expr> \)

```

```

43             | ([^--+(0-9)+) (?{
44                                     # is + and not * to avoid infinite recursion
45                                     warn "Error: expecting a number or a open parent
46                                     $warnings++;
47                                     $errors++;
48                                     }) (*COMMIT) <MATCH=factor>

```

en este caso la conducta es abandonar en el caso de que no se pueda encontrar un <factor>:

```

pl@nereida:~/Lregexpgrammars/demo/calculator$ ./calculatorwitherrmanagement.pl
Parsing infix arithmetic expressions (CTRL-D to end in unix) 2*(3+#)
Error: expecting a number or a open parenthesis, found: '#'
1 warning
1 error
Parsing infix arithmetic expressions (CTRL-D to end in unix) 2*3
postfix: 2 3 *
result: 6
Parsing infix arithmetic expressions (CTRL-D to end in unix) @
Error: expecting a number or a open parenthesis, found: '@'
1 warning
1 error
Parsing infix arithmetic expressions (CTRL-D to end in unix) Bye!

```

3.10.16. Mensajes de Warning

Sometimes, you want to detect problems, but not invalidate the entire parse as a result. For those occasions, the module provides a less stringent form of error reporting: the <warning:...> directive.

This directive is exactly the same as an <error:...> in every respect except that it does not induce a failure to match at the point it appears.

The directive is, therefore, useful for reporting non-fatal problems in a parse. For example:

```

qr{ \A          # ...Match only at start of input
    <ArithExpr>  # ...Match a valid arithmetic expression

    (?:
        # Should be at end of input...
        \s* \Z
    |
        # If not, report the fact but don't fail...
        <warning: Expected end-of-input>
        <warning: (?{ "Extra junk at index $INDEX: $CONTEXT" }>>
    )

    # Rule definitions here...
}xms;

```

Note that, because they do not induce failure, two or more <warning:...> directives can be "stacked" in sequence, as in the previous example.

3.10.17. Simplificando el AST

```

pl@nereida:~/Lregexpgrammars/demo$ cat -n exprdamian.pl
1  use strict;

```

```

2  use warnings;
3  use 5.010;
4  use Data::Dumper;
5  $Data::Dumper::Indent = 1;
6
7  my $rbb = do {
8      use Regexp::Grammars;
9
10     qr{
11         \A<expr>\z
12
13         <objrule: expr>    <MATCH=term> (?! <addop> )           # bypass
14                           | <[operands=term]> ** <[operators=addop]>
15
16         <objrule: term>    <MATCH=factor> (?! <mulop> )         # bypass
17                           | <[operands=factor]> ** <[operators=mulop]>
18
19         <objrule: factor>  <val=( [+ - ] ? \d + ( ? : \. \d * ) ? ) >
20                           | \ ( <MATCH=expr> \ )
21
22         <token: addop>    [ + - ]
23
24         <token: mulop>    [ * / ]
25
26     }x;
27 };
28
29 while (my $input = <>) {
30     chomp($input);
31     if ($input =~ m{$rbb}) {
32         my $tree = $/{expr};
33         say Dumper $tree;
34         say $tree->ceval;
35     }
36     else {
37         say("does not match");
38     }
39 }
40 }
41
42 BEGIN {
43
44     package LeftBinaryOp;
45     use strict;
46     use base qw(Class::Accessor);
47
48     LeftBinaryOp->mk_accessors(qw{operators operands});
49
50     my %f = (
51         '+' => sub { shift() + shift() },
52         '-' => sub { shift() - shift() },
53         '*' => sub { shift() * shift() },
54         '/' => sub { shift() / shift() },

```



```

55     );
56
57     sub ceval {
58         my $self = shift;
59
60         # recursively evaluate the children first
61         my @operands = map { $_->ceval } @{$self->operands};
62
63         # then combine them
64         my $s = shift @operands;
65         for (@{$self->operators}) {
66             $s = ${$_}->($s, shift @operands);
67         }
68         return $s;
69     }
70
71     package term;
72     use base qw{LeftBinaryOp};
73
74     package expr;
75     use base qw{LeftBinaryOp};
76
77     package factor;
78
79     sub ceval {
80         my $self = shift;
81
82         return $self->{val};
83     }
84
85     1;
86 }

```

Ejecuciones:

```

pl@nereida:~/Lregexpggrammars/demo$ perl5.10.1 exprdamian.pl
4-2-2

```

```

$VAR1 = bless( {
  'operands' => [
    bless( {
      '' => '4',
      'val' => '4'
    }, 'factor' ),
    bless( {
      '' => '2',
      'val' => '2'
    }, 'factor' ),
    bless( {
      '' => '2',
      'val' => '2'
    }, 'factor' )
  ],
  '' => '4-2-2',
  'operators' => [

```

```

    '-',
    '-'
```

]

```

}, 'expr' );

0
8/4/2
$VAR1 = bless( {
    'operands' => [
        bless( {
            '' => '8',
            'val' => '8'
        }, 'factor' ),
        bless( {
            '' => '4',
            'val' => '4'
        }, 'factor' ),
        bless( {
            '' => '2',
            'val' => '2'
        }, 'factor' )
    ],
    '' => '8/4/2',
    'operators' => [
        '/',
        '/'
    ]
}, 'term' );
```

```

1
3
$VAR1 = bless( {
    '' => '3',
    'val' => '3'
}, 'factor' );
```

```

3
2*(3+4)
$VAR1 = bless( {
    'operands' => [
        bless( {
            '' => '2',
            'val' => '2'
        }, 'factor' ),
        bless( {
            'operands' => [
                bless( {
                    '' => '3',
                    'val' => '3'
                }, 'factor' ),
                bless( {
                    '' => '4',
                    'val' => '4'
                }, 'factor' )
            ]
        }, 'factor' )
    ]
}, 'term' );
```

```

        }, 'factor' )
    ],
    '' => '3+4',
    'operators' => [
        '+'
    ]
}, 'expr' )
],
'' => '2*(3+4)',
'operators' => [
    '*'
]
}, 'term' );

```

14

3.10.18. Reciclando una Regexp::Grammar

Ejecución

El siguiente programa `calculator.pl` recibe como entrada una expresión en infijo.

La ejecución consta de dos bucles. En la primera parte se inyecta a la jerarquía de clases de los AST generados para las expresiones en infijo una semántica que permite evaluar la expresión:

```

58 require EvalCalc;
59
60 test_calc(
61     'Evaluating infix arithmetic expressions (CTRL-D to end in unix) ',
62     sub { print &Data::Dumper::Dumper(shift()) },
63 );

```

En esta primera parte mostraremos además el AST construido para la expresión infija de entrada.

```

pl@nereida:~/Lregexgrammars/demo$ ./calculator.pl
Evaluating infix arithmetic expressions (CTRL-D to end in unix)
8-4-2
$VAR1 = bless( {
  'operands' => [
    bless( {
      'operands' => [
        bless( {
          'operands' => [
            bless( {
              'operands' => [
                bless( { '' => '8', 'val' => '8' }, 'factor' )
              ],
              '' => '8'
            }, 'factorial' )
          ],
          '' => '8'
        }, 'power' )
      ],
      '' => '8'
    }, 'uneg' )
  ],

```

```

],
  '' => '8',
}, 'term' ),
bless( {
  'operands' => [
    bless( {
      'operands' => [
        bless( {
          'operands' => [
            bless( {
              'operands' => [
                bless( { '' => '4', 'val' => '4' }, 'factor' )
              ],
              '' => '4',
            }, 'factorial' )
          ],
          '' => '4',
        }, 'power' )
      ],
      '' => '4',
    }, 'neg' )
  ],
  '' => '4',
}, 'term' ),
bless( {
  'operands' => [
    bless( {
      'operands' => [
        bless( {
          'operands' => [
            bless( { '' => '2', 'val' => '2' }, 'factor' )
          ],
          '' => '2',
        }, 'factorial' )
      ],
      '' => '2',
    }, 'power' )
  ],
  '' => '2',
}, 'neg' )
],
  '' => '2',
}, 'term' )
],
'' => '8-4-2',
'operators' => [
  '-',
  '-',
]
}, 'expr' );
2

```

Observamos que la asociatividad es la correcta. El 2 final es el resultado de la evaluación de 8-4-2.

La estructura del árbol se corresponde con la de la gramática:

```
8 my $rbb = do {
9     use Regexp::Grammars;
10
11     qr{
12         \A<expr>\z
13
14         <objrule: expr>      <[operands=term]> ** <[operators=addop]>
15
16         <objrule: term>      <[operands=uneg]> ** <[operators=mulop]>
17
18         <objrule: uneg>      <[operators=minus]>* <[operands=power]>
19
20         <objrule: power>     <[operands=factorial]> ** <[operators=powerop]>
21
22         <objrule: factorial> <[operands=factor]> <[operators=(!)]>*
23
24         <objrule: factor>    <val=( [+ - ] ? \d + ( ? : \. \d * ) ? ) >
25                             | \ ( <MATCH=expr> \ )
26
27         <token: addop>       [ + - ]
28
29         <token: mulop>       [ * / ]
30
31         <token: powerop>     \ * \ * | \ ^
32
33         <token: minus>       - <MATCH=( ? { 'NEG' } ) >
34
35     }x;
36 };
```

Ahora, en una segunda parte sobreescribimos los métodos `sem` que describen la semántica para producir una traducción de infijo a postfijo:

```
66 require PostfixCalc;
67 test_calc('Translating expressions to postfix (CTRL-D to end in unix) ');
```

Ahora al proporcionar la entrada 6--3! obtenemos:

```
Translating expressions to postfix (CTRL-D to end in unix)
6--3!
6 3 ! ~ -
```

Aquí `~` es el operador de negación unaria y `!` es el operador factorial.

Estructura de la aplicación

Estos son los ficheros que integran la aplicación:

```
pl@nereida:~/Lregexgrammars/demo/calculator$ tree
.
|-- EvalCalc.pm          # Soporte para la evaluación de la expresión: sem
|-- Operator.pm          # Soporte a las clases nodo: recorridos
|-- PostfixCalc.pm       # Soporte para la traducción a postfijo: sem
'-- calculator.pl        # programa principal
```

Programa principal

En el programa principal definimos la gramática y escribimos una subrutina `test_calc` que realiza el parsing.

```
pl@nereida:~/Lregexgrammars/demo/calculator$ cat -n calculator.pl
1  #!/usr/bin/env perl5.10.1
2  use strict;
3  use warnings;
4  use 5.010;
5  use Data::Dumper;
6  $Data::Dumper::Indent = 1;
7
8  my $rbb = do {
9      use Regexp::Grammars;
10
11      qr{
12          \A<expr>\z
13
14          <objrule: expr>      <[operands=term]> ** <[operators=addop]>
15
16          <objrule: term>      <[operands=uneg]> ** <[operators=mulop]>
17
18          <objrule: uneg>      <[operators=minus]>* <[operands=power]>
19
20          <objrule: power>     <[operands=factorial]> ** <[operators=powerop]>
21
22          <objrule: factorial> <[operands=factor]> <[operators=(!)]>*
23
24          <objrule: factor>    <val=( [+ - ] ? \d + ( ? : \. \d * ) ? ) >
25                               | \ ( <MATCH=expr> \ )
26
27          <token: addop>       [ + - ]
28
29          <token: mulop>       [ * / ]
30
31          <token: powerop>     \ * \ * | \ ^
32
33          <token: minus>       - <MATCH=( ? { 'NEG' } ) >
34
35      }x;
36  };
37
38  sub test_calc {
39      my $prompt = shift;
40      my $handler = shift;
41
42      say $prompt;
43      while (my $input = <>) {
44          chomp($input);
45          if ($input =~ m{$rbb}) {
46              my $tree = $/{expr};
47              $handler->($tree) if $handler;
48
49              say $tree->ceval;
```



```

31     my $self = shift;
32
33     return () unless exists $self->{operators};
34     return @{$self->{operators}};
35 }
36
37 sub sem {
38     confess "not defined sem";
39 }
40
41 sub make_sem {
42     my $class = shift;
43     my %semdesc = @_;
44
45     for my $class (keys %semdesc) {
46         my %sem = %{ $semdesc{$class} };
47
48         # Install 'sem' method in $class
49         no strict 'refs';
50         no warnings 'redefine';
51         *{$class."::sem"} = sub {
52             my ($self, $op) = @_;
53             $sem{$op}
54         };
55     }
56 }
57
58 package LeftBinaryOp;
59 use base qw{Operator};
60
61 sub ceval {
62     my $self = shift;
63
64     # recursively evaluate the children first
65     my @operands = map { $_->ceval } $self->Operands;
66
67     # then combine them
68     my $s = shift @operands;
69     for ($self->Operators) {
70         $s = $self->sem($_)->($s, shift @operands);
71     }
72     return $s;
73 }
74
75 package RightBinaryOp;
76 use base qw{Operator};
77
78 sub ceval {
79     my $self = shift;
80
81     # recursively evaluate the children first
82     my @operands = map { $_->ceval } $self->Operands;
83

```



```

84     # then combine them
85     my $s = pop @operands;
86     for (reverse $self->Operators) {
87         $s = $self->sem($_)->(pop @operands, $s);
88     }
89     return $s;
90 }
91
92 package PreUnaryOp;
93 use base qw{Operator};
94
95 sub ceval {
96     my $self = shift;
97
98     # recursively evaluate the children first
99     my @operands = map { $_->ceval } $self->Operands;
100
101     # then combine them
102     my $s = shift @operands;
103     for (reverse $self->Operators) {
104         $s = $self->sem($_)->($s);
105     }
106     return $s;
107 }
108
109 package PostUnaryOp;
110 use base qw{Operator};
111
112 sub ceval {
113     my $self = shift;
114
115     # recursively evaluate the children first
116     my @operands = map { $_->ceval } $self->Operands;
117
118     # then combine them
119     my $s = shift @operands;
120     for ($self->Operators) {
121         $s = $self->sem($_)->($s);
122     }
123     return $s;
124 }
125
126 package term;
127 use base qw{LeftBinaryOp};
128
129 package expr;
130 use base qw{LeftBinaryOp};
131
132 package power;
133 use base qw{RightBinaryOp};
134
135 package uneg;
136 use base qw{PreUnaryOp};

```

```

137
138 package factorial;
139 use base qw{PostUnaryOp};
140
141 package factor;
142
143 sub ceval {
144     my $self = shift;
145
146     return $self->{val};
147 }
148
149 1;

```

Definiendo sem para la evaluación de la expresión

pl@nereida:~/Lregexpgrammars/demo/calculator\$ cat -n EvalCalc.pm

```

 1 package EvalCalc;
 2 use strict;
 3 use Carp;
 4
 5 use Operator;
 6
 7 #####
 8 sub f {
 9     $_[0]>1?$_[0]*f($_[0]-1):1;
10 }
11
12 sub fac {
13     my $n = shift;
14
15     confess "Not valid number" unless $n =~ /\d+$/;
16     f($n);
17 };
18
19 my $s = sub { shift() ** shift() };
20
21 Operator->make_sem(
22     expr => {
23         '+' => sub { shift() + shift() },
24         '-' => sub { shift() - shift() },
25     },
26     term => {
27         '*' => sub { shift() * shift() },
28         '/' => sub { shift() / shift() },
29     },
30     power => {
31         '^' => $s,
32         '**' => $s,
33     },
34     uneg => {
35         'NEG' => sub { -shift() },
36     },
37     factorial => {

```

```

38     '!' => \&fac,
39 },
40 );
41
42 1;

```

Definiendo sem para la traducción a postfijo

```

pl@nereida:~/Lregexpgrammars/demo/calculator$ cat -n PostfixCalc.pm
 1 package PostfixCalc;
 2 use strict;
 3
 4 use Operator;
 5
 6 # Modify semantics: now translate to postfix
 7 my $powers = sub { shift().' '.shift().' **' };
 8
 9 Operator->make_sem(
10     expr => {
11         '+' => sub { shift().' '.shift().' +' },
12         '-' => sub { shift().' '.shift().' -' },
13     },
14     term => {
15         '*' => sub { shift().' '.shift().' *' },
16         '/' => sub { shift().' '.shift().' /' },
17     },
18     power => {
19         '^' => $powers,
20         '**' => $powers,
21     },
22     uneg => {
23         # use ~ for unary minus
24         'NEG' => sub { shift().' ~' },
25     },
26     factorial => {
27         '!' => sub { shift().' !' },
28     },
29 );
30
31 1;

```

Ejercicio 3.10.2. ■ *Explique el significado de la primera línea del programa principal*

```

pl@nereida:~/Lregexpgrammars/demo$ cat -n calculator.pl
 1 #!/usr/bin/env perl5.10.1

```

■ *Explique el significado de \$handler en test_calc:*

```

42 sub test_calc {
43     my $prompt = shift;
44     my $handler = shift;
45
46     say $prompt;
47     while (my $input = <>) {

```

```

48     chomp($input);
49     if ($input =~ m{$rbb}) {
50         my $tree = $/{expr};
51         $handler->($tree) if $handler;
52
53         say $tree->ceval;
54
55     }
56     else {
57         say("does not match");
58     }
59 }
60 }

```

- Aísla las funciones relacionadas con la creación de semántica como `make_sem`, `fac` y las llamadas a `make_sem` en un módulo `Calculator::Semantics` aparte.
- Añade un traductor de infijo a prefijo al código presentado en esta sección. Una expresión como `2*3+4` se traducirá como `+ * 2 3 4`

3.10.19. Práctica: Calculadora con `Regexp::Grammars`

- Reforme la estructura del ejemplo para que tenga una jerarquía de desarrollo de acuerdo a los estándares de Perl. Use `h2xs` o bien `Module::Starter`. Use el espacio de nombres `Calculator`. Mueva el módulo `Operator` a `Calculator::Operator`. Lea el capítulo Módulos de los apuntes de LHP.
- Defina el conjunto de pruebas que deberá pasar su traductor. Añádalas como pruebas `TODO`. Cuando la funcionalidad a comprobar esté operativa cambie su estatus.
- Añada variables y la expresión de asignación:

```
b = a = 4*2
```

que será traducida a postfijo como:

```
4 2 * a = b =
```

El operador de asignación es asociativo a derechas. El valor devuelto por una expresión de asignación es el valor asignado.

Use un hash para implantar la relación nombre-valor en el caso de la evaluación

- Introduzca la expresión bloque:

```
c = { a = 4; b = 2*a }
```

Los bloques son listas entre llaves de expresiones separadas por punto y coma. El valor retornado por una expresión bloque es el último evaluado en el bloque.

El símbolo de arranque de la gramática (esto es, el patrón regular contra el que hay que casar) será la expresión bloque.

- Introduzca las expresiones de comparación `<`, `>`, `<=`, `>=`, `==` y `!=` con la prioridad adecuada. Tenga en cuenta que una expresión como:

`a = b+2 > c*4`

deberá entenderse como

`a = ((b+2) > (c*4))`

Esto es, se traducirá como:

`b 2 + c 4 * > a =`

- Introduzca la expresión `if ... then ... else`. La parte del `else` será opcional:

```
c = if a > 0 then { a = a -1; 2*a } else { b + 2 };
d = if a > 0 then { a = b -1; 2*b };
```

un `else` casa con el `if` mas cercano. La sentencia:

```
if (a > 0) then if (b > 0) then {5} else {6}
```

se interpreta como:

```
if (a > 0) then (if (b > 0) then {5} else {6})
```

y no como:

```
if (a > 0) then (if (b > 0) then {5}) else {6}
```

Se traducirá como:

```
      a
      0
      >
      jz endif124
      b
      0
      >
      jz else125
      5
      j endif126
:else125
      6
:endif124
:endif125
...
```

- Escriba un intérprete de la máquina orientada a pila definida en los apartados anteriores. El código generado debería poder ejecutarse correctamente en el intérprete.

Capítulo 4

La Estructura de los Compiladores: Una Introducción

Este capítulo tiene por objeto darte una visión global de la estructura de un compilador e introducirte en las técnicas básicas de la construcción de compiladores usando Perl.

Puesto que no todos los alumnos que se incorporan en este capítulo han leído los anteriores y no necesariamente conocen Perl, en la sección 4.1 comenzamos haciendo un breve repaso a como construir un módulo en Perl. Si quieres tener un conocimiento mas profundo lee el capítulo sobre módulos en [?].

La sección 4.2 describe las fases en las que -al menos conceptualmente- se divide un compilador. A continuación la sección 4.3 presenta la primera de dichas fases, el análisis léxico. En la sección 4.5 repasamos conceptos de análisis sintáctico que deberían ser familiares a cualquiera que haya seguido un curso en teoría de autómatas y lenguajes formales. Antes de comenzar a traducir es conveniente tener un *esquema* o estrategia de traducción para cada constructo sintáctico. La sección 4.7 introduce el concepto de esquema de traducción. La fase de análisis sintáctico consiste en la construcción del *árbol de análisis* a partir de la *secuencia de unidades léxicas*. Existen diversas estrategias para resolver esta fase. En la sección 4.6 introducimos la que posiblemente sea la mas sencilla de todas: el análisis descendente predictivo recursivo. En la sección 4.8 abordamos una estrategia para transformar ciertas gramáticas para las que dicho método no funciona.

Un analizador sintáctico implícitamente construye el árbol de análisis concreto. En muchas ocasiones resulta mas rentable trabajar con una forma simplificada (*abstracta*) del árbol que contiene la misma información que aquél. La sección 4.9 trata de la construcción de los árboles de análisis abstractos.

4.1. Las Bases

Puesto que no todos los alumnos que están interesados en esta sección tienen conocimientos previos de Perl, en esta sección comenzamos haciendo un breve repaso a como construir un módulo en Perl y al mismo tiempo repasamos las características usadas del lenguaje. Si quieres tener un conocimiento mas profundo de como construir un módulo, lee el capítulo sobre módulos en [?].

Version

El comportamiento de Perl puede variar ligeramente si la versión que tenemos instalada es antigua. Para ver la versión de Perl podemos hacer.

```
lhp@nereida:~/Lperl/src/topdown/PL0506$ perl -v
```

```
This is perl, v5.8.4 built for i386-linux-thread-multi
```

```
Copyright 1987-2004, Larry Wall
```

Perl may be copied only under the terms of either the Artistic License or the GNU General Public License, which may be found in the Perl 5 source kit.

Complete documentation for Perl, including FAQ lists, should be found on this system using 'man perl' or 'perldoc perl'. If you have access to the Internet, point your browser at <http://www.perl.com/>, the Perl Home Page.

h2xs

En primer lugar, construimos la estructura para nuestro proyecto de mini-lenguaje. La mejor forma de comenzar a escribir un módulo es usando la herramienta Perl **h2xs**. Supongamos que queremos construir un módulo `PL::Tutu`. Los nombres de los módulos siguen un esquema de identificadores separados por una pareja de `:`. Para saber más sobre el esquema de nombres de los módulos y la forma en la que estos se asigna a ficheros del sistema operativo, lea la sección sobre introducción a los módulos [?].

```
lhp@nereida:~/Lperl/src/topdown/PL0506$ h2xs -XA -n PL::Tutu
Defaulting to backwards compatibility with perl 5.8.4
If you intend this module to be compatible with earlier perl versions, please
specify a minimum perl version with the -b option.
```

```
Writing PL-Tutu/lib/PL/Tutu.pm
Writing PL-Tutu/Makefile.PL
Writing PL-Tutu/README
Writing PL-Tutu/t/PL-Tutu.t
Writing PL-Tutu/Changes
Writing PL-Tutu/MANIFEST
```

La herramienta **h2xs** fué concebida para ayudar en la transformación de ficheros de cabecera de C en código Perl. La opción **-X** hace que se omita la creación de subrutinas externas (XS) La opción **-A** implica que el módulo no hará uso del **AutoLoader**. La opción **-n** proporciona el nombre del módulo. La llamada a **h2xs** crea la siguiente estructura de directorios y ficheros:

```
lhp@nereida:~/Lperl/src/topdown/PL0506$ tree
.
|-- PL-Tutu
|   |-- Changes
|   |-- MANIFEST
|   |-- Makefile.PL
|   |-- README
|   |-- lib
|   |   |-- PL
|   |   |-- Tutu.pm
|-- t
|   |-- PL-Tutu.t
```

4 directories, 6 files

Generación del Makefile

Después de esto tenemos un módulo "funcional" que no hace nada. Lo podríamos instalar como si lo hubieramos descargado desde CPAN (Véase [?]).

Después cambiamos al directorio `PL-Tutu/` y hacemos `perl Makefile.PL`.

```
lhp@nereida:~/Lperl/src/topdown/PL0506$ cd PL-Tutu/
lhp@nereida:~/Lperl/src/topdown/PL0506/PL-Tutu$ perl Makefile.PL
```

```
Checking if your kit is complete...
Looks good
Writing Makefile for PL::Tutu
```

Esto crea el fichero Makefile necesario para actualizar nuestra aplicación. Para saber más sobre `perl Makefile.PL` lea [?].

Documentación

Pasamos ahora a trabajar en el módulo. Primero escribimos la parte relativa a la documentación. Para ello editamos `Tutu.pm`:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/PL-Tutu/lib/PL$ pwd
/home/lhp/Lperl/src/topdown/PL0506/PL-Tutu/lib/PL
lhp@nereida:~/Lperl/src/topdown/PL0506/PL-Tutu/lib/PL$ ls -l
total 4
-rw-r--r--  1 lhp lhp 2343 2005-09-28 11:16 Tutu.pm
```

y al final del mismo insertamos la documentación. Para saber más sobre el lenguajes de marcas de Perl (*pod* por *plain old documentation*) lea [?]. En este caso escribimos:

```
1;
__END__
```

```
=head1 NOMBRE
```

```
PL::Tutu - Compilador para un lenguaje sencillo denominado
          "Tutu" que usaremos en la asignatura PL
```

```
=head1 SINOPSIS
```

```
    use PL::Tutu;
```

```
    La subrutina PL::Tutu::compiler recibe dos argumentos: el
    nombre del fichero de entrada (fuente.tutu) y el nombre del fichero de
    salida (código ensamblador para una especie de P-máquina).
```

```
=head1 DESCRIPCIÓN
```

Este módulo tiene dos objetivos: aprender a hacer un pequeño compilador y aprender a programar modularmente en Perl, usando un buen número de los recursos que este lenguaje ofrece.

El siguiente es un ejemplo de código fuente `tutu`:

```
int a,b;
string c;
a = 2+3;
b = 3*4;
c = "hola";
p c;
c = "mundo";
p c;
p 9+2;
p a+1;
p b+1
```


supuesto que está guardado en el fichero "test2.tutu", podemos escribir un programa Perl "main.pl" para compilarlo:

```
$ cat main.pl
#!/usr/bin/perl -w -I..
#use PL::Tutu;
use Tutu;

PL::Tutu::compiler(@ARGV);
```

al ejecutar "main.pl":

```
$ ./main.pl test2.tutu test2.ok
```

obtenemos el fichero "test2.ok" con el ensamblador:

```
$ cat test2.ok
DATA holamundo
PUSH 5
PUSHADDR 0
STORE_INT
PUSH 12
PUSHADDR 1
STORE_INT
PUSHSTR 0 4
PUSHADDR 2
STORE_STRING
LOAD_STRING 2
PRINT_STR
PUSHSTR 4 5
PUSHADDR 2
STORE_STRING
LOAD_STRING 2
PRINT_STR
PUSH 11
PRINT_INT
LOAD 0
INC
PRINT_INT
LOAD 1
INC
PRINT_INT
```

Para mas información consulta la página de la asignatura.
¡Buena suerte!

=head2 EXPORT

No se exporta nada al espacio de nombres del cliente.

=head1 AUTOR

Casiano Rodríguez León, E<lt>casiano@ull.esE<gt>

=head1 VÉASE TAMBIÉN

L<perl>.

=cut

La documentación puede ser mostrada utilizando el comando `perldoc`.

Veáse la figura 4.1.

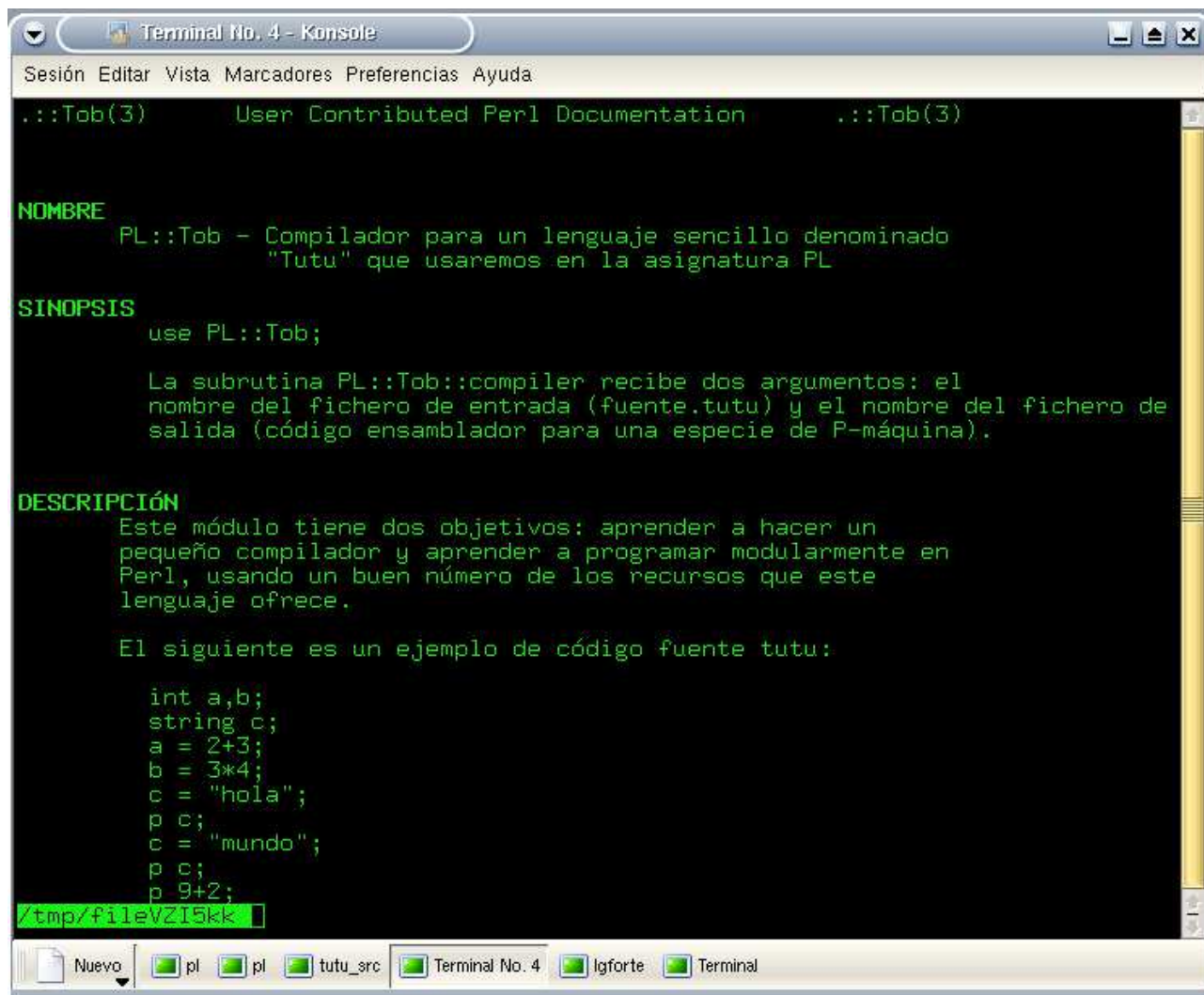


Figura 4.1: El resultado de usar `perldoc Tutu`

La forma en la que se controla la calidad de un módulo es mediante el desarrollo de pruebas. Las pruebas son programas perl que se sitúan en el directorio `t/` y que tienen la extensión `.t`. Para ejecutar las pruebas se escribe:

```
make test
```

Vease la sección [?] de los apuntes de LHP para mas detalles.

4.1.1. Repaso: Las Bases

Responda a las siguientes preguntas:

1. ¿Cómo puedo saber con que versión de Perl estoy trabajando?

2. Cuando el intérprete Perl encuentra una sentencia

```
use Este::Modulo;
```

¿Donde busca el fichero `Modulo.pm`?

3. ¿Con que opción debo usar Perl para ejecutar un programa en la línea de comandos?

4. ¿Cómo se llama el programa que me permite crear el esqueleto para una distribución de un módulo? ¿Con que opciones debo llamarlo?

5. Cuando se crea con `h2xs` el esqueleto para `PL::Tutu`: ¿En que subdirectorío queda el fichero conteniendo el esqueleto del módulo creado `Tutu.pm`?

6. ¿Cuál es la función de `MANIFEST`?

7. ¿Qué es `Makefile.PL`? ¿Cuál es su función? ¿Que significa la frase *looks good*?

8. ¿Con que comando se crea el `Makefile` para trabajar en la plataforma actual?

9. ¿Cómo se puede ver la documentación de un módulo?

10. ¿Que hacen los siguientes comandos `pod`? Repase [?] si tiene dudas.

```
=head1 cabecera
=head2 cabecera
=item texto
=over N
=back
=cut
=pod
=for X
=begin X
=end X
```

11. ¿Que secuencia de comandos conocida como *mantra de instalación* es necesario ejecutar para instalar un módulo?

12. ¿Cual es la función del directorio `t`?

13. ¿Que tipo deben tener los programas de prueba para que `make test` los reconozca como pruebas?

4.1.2. Práctica: Crear y documentar el Módulo `PL::Tutu`

Reproduzca los pasos explicados en la sección 4.1 creando el módulo `PL::Tutu` y documentándolo.

- Compruebe que la documentación se muestra correctamente
- Compruebe que puede crear una distribución haciendo `make dist`
- Compruebe que la distribución creada puede instalarse correctamente siguiendo las instrucciones en [?].

4.2. Las Fases de un Compilador

La estructura del compilador, descompuesto en fases, queda explicitada en el código de la subrutina `compile`:

Esquema del Compilador

```
1 package PL::Tutu;
2 use 5.008004;    # Versión mínima de Perl 5.84
3 use strict;      # Variables deben ser declaradas, etc.
4 use warnings;    # Enviar warnings
5 use IO::File;
6 use Carp;        # Provee alternativas a "die" and "warn"
7
8 require Exporter;
9
10 our @ISA = qw(Exporter);    # Heredamos los métodos de la clase Exporter
11 our @EXPORT = qw( compile compile_from_file); # Estas funciones serán exportadas
12 our $VERSION = '0.01';      # Variable que define la versión del módulo
13
14 our %symbol_table;          # La tabla de símbolos $symbol_table{x} contiene
15 our $data;                  # la información asociada con el objeto 'x'
16 our $target;                # tipo, dirección, etc.
17 our @tokens;                # La lista de terminales
18 our $errorflag;
19 our ($lookahead, $value);    # Token actual y su atributo
20 our $tree;                  # referencia al objeto que contiene
21 our $global_address;        # el árbol sintáctico
22
23 # Lexical analyzer
24 package Lexical::Analysis;
25 sub scanner {
26 }
27
28 package Syntax::Analysis;
29 sub parser {
30 }
31
32 package Machine::Independent::Optimization;
33 sub Optimize {
34 }
35
36 package Code::Generation;
37 sub code_generator {
38 }
39
40 package Peephole::Optimization;
41 sub transform {
42 }
43
44 package PL::Tutu;
45 sub compile {
46     my ($input) = @_; # Observe el contexto!
47     local %symbol_table = ();
48     local $data = ""; # Contiene todas las cadenas en el programa fuente
49     local $target = ""; # target code
50     local @tokens = (); # "local" salva el valor que será recuperado al finalizar
51     local $errorflag = 0; # el ámbito
52     local ($lookahead, $value) = ();
```

```

53  local $tree = undef; # Referencia al Árbol sintáctico abstracto
54  local $global_address = 0; # Usado para guardar la Última dirección ocupada
55
56  #####lexical analysis
57  &Lexical::Analysis::scanner($input);
58
59  #####syntax (and semantic) analysis
60  $tree = &Syntax::Analysis::parser;
61
62  #####machine independent optimizations
63  &Machine::Independent::Optimization::Optimize;
64
65  #####code generation
66  &Code::Generation::code_generator;
67
68  #####peephole optimization
69  &Peephole::Optimization::transform($target);
70
71  return \$target; #retornamos una referencia a $target
72 }
73
74 sub compile_from_file {
75     my ($input_name, $output_name) = @_; # Nombres de ficheros
76     my $fhi;                             # de entrada y de salida
77     my $targetref;
78
79     if (defined($input_name) and (-r $input_name)) {
80         $fhi = IO::File->new("< $input_name");
81     }
82     else { $fhi = 'STDIN'; }
83     my $input;
84     { # leer todo el fichero
85         local $/ = undef; # localizamos para evitar efectos laterales
86         $input = <$fhi>;
87     }
88     $targetref = compile($input);
89
90     #####code output
91     my $fh = defined($output_name)? IO::File->new("> $output_name") : 'STDOUT';
92     $fh->print($$targetref);
93     $fh->close;
94     1; # El Último valor evaluado es el valor retornado
95 }
96
97 1; # El 1 indica que la fase de carga termina con éxito
98 # Sigue la documentación ...

```

Añadiendo Ejecutables

Vamos a añadir un *script* que use el módulo `PL::Tutu` para así poder ejecutar nuestro compilador:

```

lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/$ mkdir scripts
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ cd scripts/

```

A continuación creamos dos versiones del compilador `tutu.pl` y `tutu` y un programa de prueba `test01.tutu`:

```
... # despues de crear los ficheros
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ ls
test01.tutu tutu tutu.pl
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ cat tutu.pl
#!/usr/bin/perl -w -I../lib/
use PL::Tutu;

PL::Tutu::compile_from_file(@ARGV);
```

Búsqueda de Librerías en Tiempo de Desarrollo

El programa tutu ilustra otra forma de conseguir que el intérprete Perl busque por la librería que está siendo desarrollada, mediante el uso de `use lib`:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ cat tutu
#!/usr/bin/perl -w
use lib ('../lib');
use PL::Tutu;

&PL::Tutu::compile_from_file(@ARGV);
```

Una tercera forma (la que recomiendo):

```
$ export PERL5LIB=~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/lib
$ perl -MPL::Tutu -e 'PL::Tutu::compile_from_file("test01.tutu")'
```

Añadiendo Los Ejecutables al MANIFEST

Ahora tenemos que añadir estos ficheros en MANIFEST para que formen parte del proyecto. En vez de eso lo que podemos hacer es crear un fichero MANIFEST.SKIP:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ cat MANIFEST.SKIP
\.$
^\.cvsignore$
/\.cvsignore$
\.cvsignore$
CVS/[~/]+$
\.svn\b
^Makefile$
/Makefile$
^blib/
\.swp$
\.bak$
\.pdf$
\.ps$
\.sal$
pm_to_blib
\.pdf$
\.tar.gz$
\.tgz$
^META.yml$
```

Ahora al hacer

```
make manifest
```

se crea un fichero MANIFEST que contiene los caminos relativos de todos los ficheros en la jerarquía cuyos nombres no casan con una de las expresiones regulares en MANIFEST.SKIP.

Para saber mas sobre MANIFEST léa [?].

No recomiendo el uso de MANIFEST.SKIP. Prefiero un control manual de los ficheros que integran la aplicacion.

Indicando al Sistema de Distribución que los Ficheros son Ejecutables

Es necesario indicarle a Perl que los ficheros añadidos son ejecutables. Esto se hace mediante el parámetro EXE_FILES de WriteMakefile:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ cat Makefile.PL
use 5.008004;
use ExtUtils::MakeMaker;
# See lib/ExtUtils/MakeMaker.pm for details of how to influence
# the contents of the Makefile that is written.
WriteMakefile(
    NAME          => 'PL::Tutu',
    VERSION_FROM  => 'lib/PL/Tutu.pm', # finds $VERSION
    PREREQ_PM     => {}, # e.g., Module::Name => 1.1
    EXE_FILES     => [ 'scripts/tutu.pl', 'scripts/tutu' ],
    ($] >= 5.005 ? ## Add these new keywords supported since 5.005
        (ABSTRACT_FROM => 'lib/PL/Tutu.pm', # retrieve abstract from module
         AUTHOR       => 'Lenguajes y Herramientas de Programacion <lhp@>') : ()),
);
```

Perl utilizará esa información durante la fase de instalación para instalar los ejecutables en el *path* de búsqueda.

Reconstrucción de la aplicación

A continuación hay que rehacer el Makefile:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ perl Makefile.PL
Writing Makefile for PL::Tutu
```

Para crear una versión funcional hacemos make:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ make
cp scripts/tutu blib/script/tutu
/usr/bin/perl "-MExtUtils::MY" -e "MY->fixin(shift)" blib/script/tutu
cp scripts/tutu.pl blib/script/tutu.pl
/usr/bin/perl "-MExtUtils::MY" -e "MY->fixin(shift)" blib/script/tutu.pl
Manifesting blib/man3/PL::Tutu.3pm
```

Para crear el MANIFEST hacemos make manifest:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ make manifest
/usr/bin/perl "-MExtUtils::Manifest=mkmanifest" -e mkmanifest
Added to MANIFEST: scripts/tutu
```

Comprobemos que el test de prueba generado automáticamente por h2xs se pasa correctamente:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ make test
PERL_DL_NONLAZY=1 /usr/bin/perl "-MExtUtils::Command::MM" "-e" "test_harness(0, 'blib/lib', 'b
t/PL-Tutu....ok
All tests successful.
Files=1, Tests=1, 0 wallclock secs ( 0.08 cusr + 0.00 csys = 0.08 CPU)
```

Ejecución

Podemos ahora ejecutar los guiones:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ cd scripts/
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ ls -l
total 12
-rw-r--r--  1 lhp lhp 15 2005-09-29 12:56 test01.tutu
-rwxr-xr-x  1 lhp lhp 92 2005-09-29 13:29 tutu
-rwxr-xr-x  1 lhp lhp 80 2005-09-29 12:58 tutu.pl
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ tutu test01.tutu test01.sal
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ ls -l
total 12
-rw-r--r--  1 lhp lhp  0 2005-09-29 13:53 test01.sal
-rw-r--r--  1 lhp lhp 15 2005-09-29 12:56 test01.tutu
-rwxr-xr-x  1 lhp lhp 92 2005-09-29 13:29 tutu
-rwxr-xr-x  1 lhp lhp 80 2005-09-29 12:58 tutu.pl
```

Veamos los contenidos del programa fuente `test01.tutu` que usaremos para hacer una prueba:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu/scripts$ cat test01.tutu
int a,b;
a = 4
```

Construcción de una Distribución

Para hacer una distribución instalable hacemos `make dist`:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ make dist
rm -rf PL-Tutu-0.01
/usr/bin/perl "-MExtUtils::Manifest=manicopy,maniread" \
    -e "manicopy(maniread(), 'PL-Tutu-0.01', 'best');"
mkdir PL-Tutu-0.01
mkdir PL-Tutu-0.01/scripts
mkdir PL-Tutu-0.01/lib
mkdir PL-Tutu-0.01/lib/PL
mkdir PL-Tutu-0.01/t
tar cvf PL-Tutu-0.01.tar PL-Tutu-0.01
PL-Tutu-0.01/
PL-Tutu-0.01/scripts/
PL-Tutu-0.01/scripts/test01.tutu
PL-Tutu-0.01/scripts/tutu
PL-Tutu-0.01/scripts/tutu.pl
PL-Tutu-0.01/META.yml
PL-Tutu-0.01/Changes
PL-Tutu-0.01/MANIFEST
PL-Tutu-0.01/lib/
PL-Tutu-0.01/lib/PL/
PL-Tutu-0.01/lib/PL/Tutu.pm
PL-Tutu-0.01/MANIFEST.SKIP
PL-Tutu-0.01/t/
PL-Tutu-0.01/t/PL-Tutu.t
PL-Tutu-0.01/Makefile.PL
PL-Tutu-0.01/README
rm -rf PL-Tutu-0.01
gzip --best PL-Tutu-0.01.tar
```


Después de esto tenemos en el directorio de trabajo el fichero `PL-Tutu-0.01.tar.gz` con la distribución:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/02fases/PL-Tutu$ ls -ltr
total 72
drwxr-xr-x  2 lhp lhp  4096 2005-09-29 12:01 t
-rw-r--r--  1 lhp lhp  1196 2005-09-29 12:01 README
drwxr-xr-x  3 lhp lhp  4096 2005-09-29 12:01 lib
-rw-r--r--  1 lhp lhp   152 2005-09-29 12:01 Changes
-rw-r--r--  1 lhp lhp   167 2005-09-29 13:23 MANIFEST.SKIP
-rw-r--r--  1 lhp lhp     0 2005-09-29 13:23 pm_to_blib
drwxr-xr-x  6 lhp lhp  4096 2005-09-29 13:23 blib
-rw-r--r--  1 lhp lhp   113 2005-09-29 13:23 MANIFEST.bak
drwxr-xr-x  2 lhp lhp  4096 2005-09-29 13:29 scripts
-rw-r--r--  1 lhp lhp   616 2005-09-29 13:49 Makefile.PL
-rw-r--r--  1 lhp lhp 20509 2005-09-29 13:51 Makefile
-rw-r--r--  1 lhp lhp  3654 2005-09-29 16:34 PL-Tutu-0.01.tar.gz
-rw-r--r--  1 lhp lhp   298 2005-09-29 16:34 META.yml
-rw-r--r--  1 lhp lhp   205 2005-09-29 16:34 MANIFEST
```

4.2.1. Repaso: Fases de un Compilador

1. ¿Que hace la declaración `package nombredepaquete`?
2. ¿Cual es la función de la declaración `use 5.008004`?
3. ¿Cuál es la función de la declaración `use strict`?
4. ¿Cuál es la función de la declaración `use warnings`?
5. ¿Que diferencia hay entre `use warnings` y `perl -w`?
6. ¿Cuál es la función de la declaración `use Carp`? ¿Que diferencia hay entre `croak` y `die`?
7. ¿Qué hace la declaración `our`?
8. ¿Qué es una variable de paquete?
9. ¿Cuál es el nombre completo de una variable de paquete?
10. ¿En que variable especial se situán los argumentos pasados a una subrutina?
11. ¿Que hace la declaración `local`?
12. ¿Cómo se declara una variable léxica?
13. ¿Cuál es el prefijo para los hashes?
14. ¿Cómo se hace referencia a un elemento de un hash `%h` de clave `k`?
15. ¿Cómo se hace referencia a un elemento de un array `@a` de índice `i`? ¿Que lugar ocupa ese elemento en el array?
16. ¿Cuál es el significado de `undef`?
17. ¿Cuál es el prefijo para las subrutinas?
18. Señale la diferencia entre

```
my ($input) = @_;
```

y

```
my $input = @_;
```

Repase [?].

19. Toda referencia es un escalar: ¿Cierto o falso?
20. Toda referencia es verdadera ¿Cierto o falso?
21. ¿Que diferencia hay entre `use` y `require`? ¿La línea `require Exporter` se ejecuta en tiempo de compilación o en tiempo de ejecución?
22. ¿Que hace la línea `our @ISA = qw(Exporter)?`. Repase [?].
23. ¿Que hace la línea `our @EXPORT = qw(compile compile_from_file)?`
24. ¿Que diferencia hay entre `EXPORT` y `EXPORT_OK`? Repase [?].
25. ¿Que hace la línea `our $VERSION = '0.01'`?
26. ¿Que valor tiene una variable no inicializada? ¿y si es un array?
27. ¿Que es un array anónimo? (Repase [?])
28. ¿Que es un hash anónimo? (Repase [?])
29. ¿Que hace el operador `=>`? Repase [?].
30. ¿En que lugar se dejan los ejecutables asociados con una distribución? ¿Cómo se informa a Perl que se trata de ejecutables?
31. ¿Cuál es la función de `MANIFEST.SKIP`? ¿Que hace `make manifest`?
32. ¿Que hace la opción `-I`? ¿Porqué la primera línea de `tutu.pl` comienza:
`#!/usr/bin/perl -w -I../lib/`
33. ¿Cómo puedo saber lo que hace el módulo `lib`? ¿Que hace la línea `use lib ('../lib')` en el programa `tutu`?
34. ¿Que contiene la variable `PERL5LIB`?
35. ¿Cómo se crea una distribución?
36. ¿Que devuelve `-r $input_name` en la línea 79? Repase [?].
37. ¿Cuál es la función de la variable mágica `$/`? ¿Que se leerá en la línea 86

```
85    local $/ = undef;  
86    my $input = <$fhi>;
```

38. ¿Que hace el operador `\`? ¿Que relación hay entre `\$target` y `$target`?
39. Si `$targetref` es una referencia a la cadena que va a contener el código objeto, ¿Cómo se denota a la cadena referenciada por `$targetref`? Explique la línea

```
92    $fh->print($$targetref);
```

4.2.2. Práctica: Fases de un Compilador

Reproduzca los pasos explicados en la sección 4.2 extendiendo el módulo `PL::Tutu` con las funciones de compilación y los correspondientes guiones de compilación.

Mejore el script `tutu` para que acepte opciones desde la línea de comandos. Debera soportar al menos las siguientes opciones:

- `--usage`
Muestra de forma concisa el comando de uso
- `--help`
Un resumen de cada opción disponible
- `--version`
Muestra la versión del programa
- `--man`
Muestra la documentación

Use para ello el módulo `Getopt::Long`. Este módulo provee la función `GetOptions` la cual se atiene a los estándares de especificación de opciones en la línea de comandos POSIX y GNU. Esta función soporta el uso del guión doble `--` y el simple así como admitir el prefijo mas corto que deshace la ambigüedad entre las diferentes opciones.

La llamada a `GetOptions` analiza la línea de comandos en `ARGV` inicializa la variable asociada de manera adecuada. Retorna un valor verdadero si la línea de comandos pudo ser procesada con En caso contrario emitirá un mensaje de error y devolverá falso. Recuerde hacer `perldoc Getopt::Long` para obtener información mas detallada

El siguiente ejemplo ilustra el uso de `Getopt::Long`. Se hace uso también del módulo (función `pod2usage` en la línea 63) `Pod::Usage` el cual permite la documentación empotrada.

```
nereida:~/LEyapp/examples> cat -n treereg
 1  #!/usr/bin/perl -w
 2  use strict;
 3  use Parse::Eyapp::YATW;
 4  use Parse::Eyapp::Node;
 5  use Parse::Eyapp::Treeregexp;
 6  use Carp;
 7  use Getopt::Long;
 8  use Pod::Usage;
 9
10  my $infile;
11  my $outfile;
12  my $packagename;
13  my $prefix = '';
14  my $syntax = 1;
15  my $numbers = 1;
16  my $severity = 0; # 0 = Don't check arity. 1 = Check arity.
17                  # 2 = Check arity and give a warning 3 = ... and croak
18  GetOptions(
19    'in=s'      => \$infile,
20    'out=s'     => \$outfile,
21    'mod=s'     => \$packagename,
22    'prefix=s'  => \$prefix,
23    'severity=i'=> \$severity,
24    'syntax!'   => \$syntax,
```

```

25     'numbers!'    => \$numbers,
26     'version'    => \&version,
27     'usage'      => \&usage,
28     'help'       => \&man,
29 ) or croak usage();
30
31 # If an argument remains is the inputfile
32 ($infile) = @ARGV unless defined($infile);
33 die usage() unless defined($infile);
34
35 my $treeparser = Parse::Eyapp::Treeregexp->new(
36     INFILE    => $infile,
37     OUTFILE   => $outfile,
38     PACKAGE   => $packagename,
39     PREFIX    => $prefix,
40     SYNTAX    => $syntax,
41     NUMBERS   => $numbers,
42     SEVERITY  => $severity
43 );
44
45 $treeparser->generate();
46
47 sub version {
48     print "Version $Parse::Eyapp::Treeregexp::VERSION\n";
49     exit;
50 }
51
52 sub usage {
53     print <<"END_ERR";
54     Supply the name of a file containing a tree grammar (.trg)
55     Usage is:
56     treereg [-m packagename] [[no]syntax] [[no]numbers] [-severity 0|1|2|3] \
57         [-p treeprefix] [-o outputfile] -i filename[.trg]
58     END_ERR
59     exit;
60 }
61
62 sub man {
63     pod2usage(
64         -exitval => 1,
65         -verbose => 2
66     );
67 }
68 __END__
69
70 =head1 SYNOPSIS
71
72     treereg [-m packagename] [[no]syntax] [[no]numbers] [-severity 0|1|2|3] \
73         [-p treeprefix] [-o outputfile] -i filename[.trg]
74     treereg [-m packagename] [[no]syntax] [[no]numbers] [-severity 0|1|2|3] \
75         [-p treeprefix] [-o outputfile] filename[.trg]
76 ... # Follows the documentation bla, bla, bla

```

Ahora podemos ejecutar el guión de múltiples formas:

```

nereida:~/LEyapp/examples> treereg -nos -nonu -se 3 -m Tutu Foldonly1.trg
nereida:~/LEyapp/examples> treereg -nos -nonu -s 3 -m Tutu Foldonly1.trg
Option s is ambiguous (severity, syntax)
nereida:~/LEyapp/examples> treereg -nos -bla -nonu -m Tutu Foldonly1.trg
Unknown option: bla
nereida:~/LEyapp/examples>

```

La librería estandar de Perl incluye el módulo `Getopt::Long`. No es el caso de `Pod::Usage`. Descarga el módulo e instalo en un directorio local en el que tenga permisos. Si es preciso repase las secciones [?] y [?] de los apuntes de introducción a Perl.

4.3. Análisis Léxico

Comenzaremos con la parte mas sencilla del compilador: el analizador léxico. Habitualmente el término “análisis léxico” se refiere al tratamiento de la entrada que produce como salida la lista de *tokens*. Un *token* hace alusión a las unidades mas simples que tiene significado. Habitualmente un *token* o lexema queda descrito por una expresión regular. Léxico viene del griego *lexis*, que significa “palabra”. Perl es, sobra decirlo, una herramienta eficaz para encontrar en que lugar de la cadena se produce un emparejamiento. Sin embargo, en el análisis léxico, el problema es encontrar la subcadena a partir de la última posición en la que se produjo un emparejamiento y que es aceptada por una de las expresiones regulares que definen los lexemas del lenguaje dado.

La estructura general del analizador léxico consiste en un bucle en el que se va recorriendo la entrada, buscando por un emparejamiento con uno de los patrones/lexemas especificados y, cuando se encuentra, se retorna esa información al analizador sintáctico. Como no tenemos escrito el analizador sintáctico simplemente iremos añadiendo los terminales al final de una lista.

Una iteración del bucle tiene la forma de una secuencia de condicionales en las que se va comprobando si la entrada casa con cada una de las expresiones regulares que definen los terminales del lenguaje. Las condiciones tienen un aspecto similar a este:

```

...
if (m{\G\s*(\d+)}gc) {
    push @tokens, 'NUM', $1;
}
elsif (m{\G\s*([a-z_]\w*)\b}igc) {
    push @tokens, 'ID', $1;
}
...

```

Una expresión como `m{\G\s*(\d+)}gc` es una expresión regular. Es conveniente que en este punto repase la introducción a las expresiones regulares en [?].

El Operador de Binding

Nótese que, puesto que no se menciona sobre que variable se hace el *binding* (no se usa ninguno de los operadores de *binding* `=~` y `!~`): se entiende que es sobre la variable por defecto `$_` sobre la que se hace el matching.

Casando a partir del Último Emparejamiento

El ancla `\G` casa con el punto en la cadena en el que terminó el último emparejamiento.

La expresión regular describiendo el patrón de interés se pone entre paréntesis para usar la estrategia de los paréntesis con memoria (véanse 3.1.4 y 3.1.1).

Las opciones `c` y `g` son especialmente útiles para la construcción del analizador.

La opción g

Como lo usamos en un contexto escalar, la opción `g` itera sobre la cadena, devolviendo cierto cada vez que casa, y falso cuando deja de casar. Se puede averiguar la posición del emparejamiento utilizando la función `pos`. (véase la sección 3.1.6 para mas información sobre la opción `g`).

La opción c

La opción `/c` afecta a las operaciones de emparejamiento con `/g` en un contexto escalar. Normalmente, cuando una búsqueda global escalar tiene lugar y no ocurre casamiento, la posición de comienzo de búsqueda es reestablecida al comienzo de la cadena. La opción `/c` hace que la posición inicial de emparejamiento permanezca donde la dejó el último emparejamiento con éxito y no se vaya al comienzo. Al combinar esto con el ancla `\G`, la cuál casa con el final del último emparejamiento, obtenemos que la combinación

`m{\G\s*(...)}gc`

logra el efecto deseado: Si la primera expresión regular en la cadena `elsif` fracasa, la posición de búsqueda no es inicializada de nuevo gracias a la opción `c` y el ancla `\G` sigue recordando donde terminó el ultimo casamiento.

La opción i

Por último, la opción `i` permite ignorar el tipo de letra (mayúsculas o minúsculas).

Repase la sección 3.1.6 para ver algunas de las opciones mas usadas.

Código del Analizador Léxico

Este es el código completo de la subrutina `scanner` que se encarga del análisis léxico:

```
1 package Lexical::Analysis;
2 sub scanner {
3     local $_ = shift;
4     { # Con el redo del final hacemos un bucle "infinito"
5         if (m{\G\s*(\d+)}gc) {
6             push @tokens, 'NUM', $1;
7         }
8         elsif (m{\G\s*int\b}igc) {
9             push @tokens, 'INT', 'INT';
10        }
11        elsif (m{\G\s*string\b}igc) {
12            push @tokens, 'STRING', 'STRING';
13        }
14        elsif (m{\G\s*p\b}igc) {
15            push @tokens, 'P', 'P'; # P para imprimir
16        }
17        elsif (m{\G\s*([a-z_]\w*)\b}igc) {
18            push @tokens, 'ID', $1;
19        }
20        elsif (m{\G\s*"([^"]*)" }igc) {
21            push @tokens, 'STR', $1;
22        }
23        elsif (m{\G\s*([+*()=;,]}gc) {
24            push @tokens, 'PUN', $1;
25        }
26        elsif (m{\G\s*(\S)}gc) { # Hay un caracter "no blanco"
27            Error::fatal "Caracter invalido: $1\n";
28        }
```

```

29     else {
30         last;
31     }
32     redo;
33 }
34 }

```

Relación de corutina con el Analizador Sintáctico

Si decidieramos establecer una relación de corutina con el analizador léxico los condicionales se pueden programar siguiendo secuencias con esta estructura:

```
return ('INT', 'INT') if (m{\G\s*int\b}igc);
```

Manejo de Errores

Para completar el analizador solo quedan declarar las variables usadas y las subrutinas de manejo de errores:

```
##### global scope variables
our @tokens = ();
our $errorflag = 0;

package Error;

sub error($) {
    my $msg = shift;
    if (!$errorflag) {
        warn "Error: $msg\n";
        $errorflag = 1;
    }
}

sub fatal($) {
    my $msg = shift;
    die "Error: $msg\n";
}

```

El uso de `our` es necesario porque hemos declarado al comienzo del módulo `use strict`. El pragma `use strict` le indica al compilador Perl que debe considerar como obligatorias un conjunto de reglas de buen estilo de programación. Entre otras restricciones, el uso del pragma implica que todas las variables (no-mágicas) deben ser declaradas explícitamente (uso de `my`, `our`, etc.) La declaración `our` se describe en [?].

4.3.1. Ejercicio: La opción `g`

Explique cada una de las líneas que siguen:

```

$ perl -wde 0
main::(-e:1): 0
DB<1> $x = "ababab"
DB<2> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)
match= b pos = 2
DB<3> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)
match= b pos = 4
DB<4> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)

```

```
match= b pos = 6
DB<5> print "falso" unless $x =~ m{b}g
falso
```

4.3.2. Ejercicio: Opciones g y c en Expresiones Regulares

Explique cada una de las conductas que siguen

- ¿Porqué en la línea 18 se casa con la primera b?

```
DB<5> $x = "bbabab"
DB<6> $x =~ m{a}g; print "match= ".$&." pos = ".pos($x)
match= a pos = 3
DB<7> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)
match= b pos = 4
DB<8> $x =~ m{c}g; print "match= ".$&." pos = ".pos($x)
Use of uninitialized value in concatenation (.)
DB<18> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)
match= b pos = 1
```

- ¿Porqué en la línea 27 se casa con la última b?

```
DB<23> $x = "bbabab"
DB<24> $x =~ m{a}g; print "match= ".$&." pos = ".pos($x)
match= a pos = 3
DB<25> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)
match= b pos = 4
DB<26> $x =~ m{c}gc
DB<27> $x =~ m{b}g; print "match= ".$&." pos = ".pos($x)
match= b pos = 6
```

- ¿Porqué en la línea 5 se produce casamiento y en la línea 8 no?

```
DB<3> $x = "bcbabab"
DB<4> $x =~ m{b}gc; print "match= ".$&." pos = ".pos($x)
match= b pos = 1
DB<5> $x =~ m{a}gc; print "match= ".$&." pos = ".pos($x)
match= a pos = 4

DB<6> $x = "bcbabab"
DB<7> $x =~ m{b}gc; print "match= ".$&." pos = ".pos($x)
match= b pos = 1
DB<8> $x =~ m{\Ga}gc; print "match= ".$&." pos = ".pos($x)
Use of uninitialized value in concatenation
match= pos = 1
```

4.3.3. Ejercicio: El orden de las expresiones regulares

¿Que ocurriría en la subrutina `scanner` si el código en las líneas 17-19 que reconoce los identificadores se adelanta a la línea 8? ¿Que ocurriría con el reconocimiento de las palabras reservadas como `INT`? ¿Seguiría funcionando correctamente el analizador?

4.3.4. Ejercicio: Regexp para cadenas

En la rutina `scanner`. ¿Es legal que una cadena correspondiente al terminal `STR` contenga retornos de carro entre las comillas dobles?

4.3.5. Ejercicio: El or es vago

Explique el resultado de la siguiente sesión con el depurador:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL/Lexical$ perl -de 0
DB<1> 'bb' =~ m{b|bb}; print $&
b
DB<2> 'bb' =~ m{bb|b}; print $&
bb
```

Perl convierte la expresión regular en un NFA. A diferencia de lo que ocurre en otras herramientas, el NFA no es convertido en un DFA. El NFA es entonces simulado. ¿Que está ocurriendo en la simulación?

4.3.6. Práctica: Números de Línea, Errores, Cadenas y Comentarios

Extienda el analizador léxico para que:

- Reescriba la expresión regular para las cadenas de manera que acepte comillas dobles escapadas `\` en el interior de una cadena. Por ejemplo en `"esta \"palabra\" va entre comillas"`.

Analice esta solución ¿Es correcta?:

```
DB<1> $stringre = qr{"(\\.|[^\"])*"}
DB<2> print $& if '"esta \"palabra\" va entre comillas"' =~ $stringre
"esta \"palabra\" va entre comillas"
```

- Consuma comentarios a la Perl: cualesquiera caracteres después de una almohadilla hasta el final de la línea (`# ...`).
- Consuma comentarios no anidados a la C (`/* ... */`). Repase las secciones sobre expresiones regulares no “greedy” (p. ej. sección 3.1.1) y la sección 3.1.1. Recuerde que, en una expresión regular, la opción `/s` hace que el punto `.` empareje con un retorno de carro `\n`. Esto es, el punto “casa” con cualquier carácter.

Observe el siguiente ejemplo:

```
pl@nereida:~/src/perl/testing$ cat -n ccomments.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3
 4  sub showmatches {
 5      my ($x, $re) = @_;
 6
 7      for (my $i = 0; $x =~ /$re/gsx; $i++) {
 8          print "matching $i: $1\n";
 9          $i++;
10      }
11  }
12
13  my $x = <<'EOC';
14  if (x) {
15      /* a comment */ return x + 1; /* another comment */
16  }
17  else {
18      return x + 2; /* a last comment */
19  }
20  EOC
```

```

21
22 print "\n*****\n";
23
24 my $greedy = q{ (
25     /\* # Abrir comentario
26     .* # Consumir caracteres (greedy)
27     \*/ # Cerrar comentario
28 )
29 };
30 print "Greedy:\n";
31 showmatches($x, $greedy);
32
33 print "\n*****\n";
34
35 my $lazy = q{ (
36     /\* # Abrir comentario
37     .*? # Consumir caracteres (lazy)
38     \*/ # Cerrar comentario
39 )
40 };
41 print "Lazy:\n";
42 showmatches($x, $lazy);

```

Cuando se ejecuta produce:

```
pl@nereida:~/src/perl/testing$ ccomments.pl
```

```

*****
Greedy:
matching 0: /* a comment */ return x + 1; /* another comment */
}
else {
    return x + 2; /* a last comment */

*****
Lazy:
matching 0: /* a comment */
matching 2: /* another comment */
matching 4: /* a last comment */

```

Explique la conducta.

- Números en punto flotante (como `-1.32e-04` o `.91`). El siguiente ejemplo intenta ayudarle en la búsqueda de la solución:

```

lhp@nereida:~$ perl -de 0
DB<1> print "Si" if ('.5' =~ m{\d+})
Si
DB<2> print "Si" if ('.5' =~ m{^\d+})

DB<3> print "Si" if '0.7' =~ m{^\d+(\.\d+)?(e[+-]?\d+)?$}
Si
DB<4> print "Si" if '.7' =~ m{^\d+(\.\d+)?(e[+-]?\d+)?$}

```

```

DB<5> print "Si" if '1e2' =~ m{^\d+(\.\d+)?(e[+-]?\d+)?$}
Si
DB<6> print "Si " while 'ababa' =~ m{a}g
Si Si Si
DB<7> print "@a" if @a = 'ababa' =~ m{(a)}g
a a a

```

¿Sabría decir porque la respuesta al primer comando del depurador es `si`?

- Tenga presente el posible conflicto entre los terminales `INT` y `FLOAT`. Si la entrada contiene 3.5 el terminal debería ser (`FLOAT`, '3.5') y no (`INT`, 3), ('.', '.'), (`INT`, 5).
- En esta práctica si lo desea puede instalar y usar el módulo `Regexp::Common` mantenido por *Abigail* el cual provee expresiones regulares para las situaciones mas comunes: números, teléfonos, IP, códigos postales, listas, etc. Puede incluso usarlo para encontrar soluciones a las cuestiones planteadas en esta práctica:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK> perl -MRegexp::Common -e 'print "$RE{num}{int}\n"'
(?: (?: [+ -] ?) (?: [0123456789] +))

```

Podemos hacer uso directo del hash `%RE` directamente en las expresiones regulares aprovechando que estas interpolan las variables en su interior:

```

nereida:/tmp> cat -n prueba.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3  use Regexp::Common;
 4
 5  my $input = <>;
 6
 7  print "$&\n" if $input =~ /^$RE{num}{real}$/;
nereida:/tmp> ./prueba.pl
23.45
23.45
nereida:/tmp> ./prueba.pl
jshdf
nereida:/tmp>

```

- Para mejorar la calidad de los mensajes de error extienda el par (`terminal`, `valor`) devuelto por el `scanner` a un par (`terminal`, [`valor`, número de línea]) cuya segunda componente es un array anónimo conteniendo el valor y el número de línea en el que aparece el terminal.

El siguiente extracto de un analizador léxico muestra como hacerlo:

```

sub _Lexer {

    return('', undef) unless defined($input);

    #Skip blanks
    $input =~ m{\G((?:
        \s+          # any white space char
        | \#[^\n]*   # Perl like comments
    )+)}xsgc
    and do {
        my($blanks)=$1;

```

```

    #Maybe At EOF
    pos($input) >= length($input)
    and return('', undef);
    $tokenend += $blanks =~ tr/\n//;
};

$tokenbegin = $tokenend;

$input=~/\G(and)/gc
and return($1, [$1, $tokenbegin]);

$input=~/\G(?:[A-Za-z_][A-Za-z0-9_]*:*)([A-Za-z_][A-Za-z0-9_]+)/gc
and do {
    return('IDENT', [$1, $tokenbegin]);
};

.....

$input=~/\G{/gc
and do {
    my($level,$from,$code);

    $from=pos($input);

    $level=1;
    while($input=~/(\[{}])/gc) {
        substr($input,pos($input)-1,1) eq '\\"' #Quoted
        and next;
        $level += ($1 eq '{' ? 1 : -1)
        or last;
    }
    $level
    and _SyntaxError("Not closed open curly bracket { at $tokenbegin");
    $code = substr($input,$from,pos($input)-$from-1);
    $tokenend+= $code=~tr/\n//;
    return('CODE', [$code, $tokenbegin]);
};

#Always return something
$input=~/\G(.) /sg
and do {
    $1 eq "\n" and ++$tokenend;
    return ($1, [$1, $tokenbegin]);
};
#At EOF
return('', undef);
}

```

El operador `tr` ha sido utilizado para contar los retornos de carro (descrito en la sección 3.5). El operador, además de reemplazar devuelve el número de caracteres reemplazados o suprimidos:

```
$cuenta = $cielo =~ tr/*/*/; # cuenta el numero de estrellas en cielo
```

Para aprender soluciones alternativas consulte `perldoc -q 'substring'`.

- Mejore sus mensajes de error ahora que lleva la cuenta de los números de línea. En vez de usar las rutinas `error` y `fatal` introducidas en la sección anterior escriba una sola rutina que recibe el nivel de severidad del error (parámetro `$level` en el siguiente ejemplo) y ejecuta la acción apropiada. El código de `_SyntaxError` ilustra como hacerlo:

```
sub _SyntaxError {
    my($level,$message,$lineno)=@_;

    $message= "*".
        [ 'Warning', 'Error', 'Fatal' ]->[$level].
        "* $message, at ".
        ($lineno < 0 ? "eof" : "line $lineno")." at file $filename\n";

    $level > 1
    and die $message;

    warn $message;

    $level > 0 and ++$nberr;

    $nberr == $max_errors
    and die "*Fatal* Too many errors detected.\n"
}
```

4.4. Pruebas para el Analizador Léxico

Queremos separar/aislar las diferentes fases del compilador en diferentes módulos.

Módulo PL::Error

Para ello comenzamos creando un módulo conteniendo las rutinas de tratamiento de errores:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL$ pwd
/home/lhp/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL$ cat -n Error.pm
 1  package Error;
 2  use strict;
 3  use warnings;
 4  use Carp;
 5
 6  require Exporter;
 7
 8  our @ISA = qw(Exporter);
 9  our @EXPORT = qw( error fatal);
10  our $VERSION = '0.01';
11
12  sub error {
13      my $msg = join " ", @_;
14      if (!$PL::Tutu::errorflag) {
15          carp("Error: $msg\n");
16          $PL::Tutu::errorflag = 1;
17      }
18  }
```

```

19
20 sub fatal {
21     my $msg = join " ", @_;
22     croak("Error: $msg\n");
23 }

```

Observa como accedemos a la variable `errorflag` del paquete `PL::Tutu`. Para usar este módulo desde `PL::Tutu`, tenemos que declarar su uso:

```

lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL$ cat -n Tutu.pm | head -8
 1 package PL::Tutu;
 2
 3 use 5.008004;
 4 use strict;
 5 use warnings;
 6 use IO::File;
 7 use Carp;
 8 use PL::Error;

```

En la línea 8 hacemos `use PL::Error` y no `use Error` ya que el módulo lo hemos puesto en el directorio `PL`. No olvides hacer `make manifest` para actualizar el fichero `MANIFEST`.

Módulo `PL::Lexical::Analysis`

Supongamos que además de modularizar el grupo de rutinas de tratamiento de errores queremos hacer lo mismo con la parte del análisis léxico. Parece lógico que el fichero lo pongamos en un subdirectororio de `PL/` por lo que cambiamos el nombre del módulo a `PL::Lexical::Analysis` quedando la jerarquía de ficheros así:

```

lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL$ tree
.
|-- Error.pm
|-- Lexical
|   '-- Analysis.pm
'-- Tutu.pm

```

Por supuesto debemos modificar las correspondientes líneas en `Tutu.pm`:

```

 1 package PL::Tutu;
 2
 3 use 5.008004;
 4 use strict;
 5 use warnings;
 6 use IO::File;
 7 use Carp;
 8 use PL::Error;
 9 use PL::Lexical::Analysis;
10 ...
11
12 sub compile {
13     my ($input) = @_;
14     local %symbol_table = ();
15     local $data = ""; # Contiene todas las cadenas en el programa fuente
16     local $target = ""; # target code
17     my @tokens = ();
18     local $errorflag = 0;

```

```

19  local ($lookahead, $value) = ();
20  local $tree = undef; # abstract syntax tree
21  local $global_address = 0;
22
23
24  #####lexical analysis
25  @tokens = &PL::Lexical::Analysis::scanner($input);
26  print "@tokens\n";
27
28  ...
29
30  return \$target;
31 }

```

Observe que ahora `PL::Lexical::Analysis::scanner` devuelve ahora la lista con los terminales y que `@tokens` se ha ocultado en `compile` como una variable léxica (línea 17). En la línea 26 mostramos el contenido de la lista de terminales.

Sigue el listado del módulo conteniendo el analizador léxico. Obsérve las líneas 6, 16 y 44.

```

lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/lib/PL/Lexical$ cat -n Analysis.pm
 1  # Lexical analyzer
 2  package PL::Lexical::Analysis;
 3  use strict;
 4  use warnings;
 5  use Carp;
 6  use PL::Error;
 7
 8  require Exporter;
 9
10  our @ISA = qw(Exporter);
11  our @EXPORT = qw( scanner );
12  our $VERSION = '0.01';
13
14  sub scanner {
15      local $_ = shift;
16      my @tokens;
17
18      { # Con el redo del final hacemos un bucle "infinito"
19          if (m{\G\s*(\d+)}gc) {
20              push @tokens, 'NUM', $1;
21          }
22          ...
37          elsif (m{\G\s*([+*()=;,])}gc) {
38              push @tokens, 'PUN', $1;
39          }
40          elsif (m{\G\s*(\S)}gc) {
41              Error::fatal "Caracter invalido: $1\n";
42          }
43          else {
44              return @tokens;
45          }
46          redo;
47      }
48  }

```

El Programa Cliente

Puesto que en el paquete `PL::Lexical::Analysis` exportamos `scanner` no es necesario llamar la rutina por el nombre completo desde `compile`. Podemos simplificar la línea en la que se llama a `scanner` que queda así:

```
#####lexical analysis
@tokens = &scanner($input);
print "@tokens\n";
```

De la misma forma, dado que `PL::Tutu` exporta la función `compile_from_file`, no es necesario llamarla por su nombre completo desde el guión `tutu`. Reescribimos la línea de llamada:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/scripts$ cat tutu
#!/usr/bin/perl -w
use lib ('../lib');
use PL::Tutu;

&compile_from_file(@ARGV);
```

Actualización del MANIFEST

Como siempre que se añaden o suprimen archivos es necesario actualizar `MANIFEST`:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu$ make manifest
/usr/bin/perl "-MExtUtils::Manifest=mkmanifest" -e mkmanifest
Added to MANIFEST: lib/PL/Lexical/Analysis.pm
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu$ cat -n MANIFEST
 1 Changes
 2 lib/PL/Error.pm
 3 lib/PL/Lexical/Analysis.pm
 4 lib/PL/Tutu.pm
 5 Makefile.PL
 6 MANIFEST
 7 MANIFEST.SKIP
 8 README
 9 scripts/test01.tutu
10 scripts/tutu
11 scripts/tutu.pl
12 t/PL-Tutu.t
```

4.4.1. Comprobando el Analizador Léxico

Queremos comprobar si nuestro código funciona. ¿Cómo hacerlo?. Lo adecuado es llevar una aproximación sistemática que permita validar el código.

Principios Básicos del Desarrollo de Pruebas

En general, la filosofía aconsejable para realizar un banco de pruebas de nuestro módulo es la que se articula en la metodología denominada *Extreme Programming*, descrita en múltiples textos, en concreto en el libro de Scott [?]:

- Todas las pruebas deben automatizarse
- Todos los fallos que se detecten deberían quedar traducidos en pruebas
- La aplicación debería pasar todas las pruebas después de cualquier modificación importante y también al final del día
- El desarrollo de las pruebas debe preceder el desarrollo del código
- Todos los requerimientos deben ser expresados en forma de pruebas

La Jerarquía de Una Aplicación

Pueden haber algunas diferencias entre el esquema que se describe aquí y su versión de Perl. Lea detenidamente el capítulo Test Now, test Forever del libro de Scott [?] y el libro [?] de Ian Langworth y chromatic.

Si usas una versión de Perl posterior la 5.8.0, entonces tu versión del programa `h2xs` creará un subdirectorio `/t` en el que guardar los ficheros de prueba. Estos ficheros deberán ser programas Perl de prueba con el tipo `.t`. La utilidad `h2xs` incluso deja un programa de prueba `PL-Tutu.t` en ese directorio. La jerarquía de ficheros con la que trabajamos actualmente es:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu$ make veryclean
rm -f blib/script/tutu blib/script/tutu.pl
rm -rf ./blib Makefile.aperl ...
mv Makefile Makefile.old > /dev/null 2>&1
rm -rf blib/lib/auto/PL/Tutu blib/arch/auto/PL/Tutu
rm -rf PL-Tutu-0.01
rm -f blib/lib/PL/.Tutu.pm.swp ...
rm -f *~ *.orig */*~ */*.orig
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu$ tree
.
|-- .svn                # use siempre un sistema de control de versiones
|-- Changes             # la historia de cambios
|-- MANIFEST            # lista de ficheros que componen la distribución
|-- MANIFEST.SKIP       # regexps para determinar que ficheros no pertenecen
|-- META.yml            # YML no es XML
|-- Makefile.PL         # generador del Makefile independiente de la plataforma
|-- PL-Tutu-0.01.tar.gz
|-- README              # instrucciones de instalacion
|-- lib
|   '-- PL
|       |-- Error.pm    # rutinas de manejo de errores
|       |-- Lexical
|           |-- Analysis.pm # modulo con el analizador lexico
|           '-- Tutu.pm    # modulo principal
|-- scripts
|   |-- test01.sal      # salida del programa de prueba
|   |-- test01.tutu     # programa de prueba
|   |-- tutu            # compilador
|   '-- tutu.pl         # compilador
'-- t
    '-- 01Lexical.t     # prueba consolidada
```

Un Ejemplo de Programa de Prueba

Estos son los contenidos de nuestro primer test:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu$ cd t
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/t$ ls -l
total 4
-rw-r--r--  1 lhp lhp 767 2005-10-10 11:27 01Lexical.t
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/t$ cat -n 01Lexical.t
 1 # Before 'make install' is performed this script should be runnable with
 2 # 'make test'. After 'make install' it should work as 'perl PL-Tutu.t'
 3
 4 #####
 5
```

```

6  # change 'tests => 1' to 'tests => last_test_to_print';
7
8  use Test::More tests => 5;
9  use Test::Exception;
10
11 BEGIN { use_ok('PL::Lexical::Analysis') };
12 BEGIN { use_ok('PL::Tutu') };
13
14 #####
15
16 # Insert your test code below, the Test::More module is use()ed here so read
17 # its man page ( perldoc Test::More ) for help writing this test script.
18
19 can_ok('PL::Lexical::Analysis', 'scanner');
20
21 # Test result of call
22 my $a = 'int a,b; string c; c = "hello"; a = 4; b = a +1; p b';
23 my @tokens = scanner($a);
24 my @expected_tokens = qw{
25 INT INT
26 ID a
27 PUN ,
28 ID b
29 PUN ;
30 STRING STRING
31 ID c
32 PUN ;
33 ID c
34 PUN =
35 STR "hello"
36 PUN ;
37 ID a
38 PUN =
39 NUM 4
40 PUN ;
41 ID b
42 PUN =
43 ID a
44 PUN +
45 NUM 1
46 PUN ;
47 P P
48 ID b
49 };
50 is(@tokens, @expected_tokens, "lexical analysis");
51
52 # test a lexically erroneous program
53 $a = 'int a,b; string c[2]; c = "hello"; a = 4; b = a +1; p b';
54 throws_ok { scanner($a) } qr{Error: Carácter inválido:}, 'erroneous program';

```

El nombre del fichero de prueba debe cumplir que:

- Sea significativo del tipo de prueba
- Que los prefijos de los nombres 01, 02, ... nos garanticen el orden de ejecución

Ejecución de Las Pruebas

Ahora ejecutamos las pruebas:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu$ make test
PERL_DL_NONLAZY=1 /usr/bin/perl "-MExtUtils::Command::MM" "-e" "test_harness(0, 'blib/lib', 'b
t/01Lexical....ok 1/5Possible attempt to separate words with commas at t/01Lexical.t line 49.
t/01Lexical....ok
All tests successful.
Files=1, Tests=5, 0 wallclock secs ( 0.08 cusr + 0.00 csys = 0.08 CPU)
```

O bien usamos prove:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/t$ prove -I../lib 01Lexical.t
01Lexical....ok
All tests successful.
Files=1, Tests=2, 0 wallclock secs ( 0.03 cusr + 0.01 csys = 0.04 CPU)
```

También podemos añadir la opción verbose a prove:

```
lhp@nereida:~/Lperl/src/topdown/PL0506/03lexico/PL-Tutu/t$ prove -v 01Lexical.t
01Lexical....1..5
ok 1 - use PL::Lexical::Analysis;
ok 2 - use PL::Tutu;
ok 3 - PL::Lexical::Analysis->can('scanner')
ok 4 - lexical analysis
ok 5 - erroneous program
ok
All tests successful.
Files=1, Tests=5, 0 wallclock secs ( 0.07 cusr + 0.01 csys = 0.08 CPU)
```

Repáse [?] para un mejor conocimiento de la metodología de pruebas en Perl.

Versiones anteriores a la 5.8

En esta sección he usado la versión 5.6.1 de Perl. Creemos un subdirectorio `tutu_src/` y en él un programa de prueba `pruebalex.pl`:

```
$ pwd
/home/lhp/projects/perl/src/tmp/PL/Tutu/tutu_src
$ cat pruebalex.pl
#!/usr/bin/perl -w -I..
#use PL::Tutu;
use Tutu;

my $a = 'int a,b; string c; c = "hello"; a = 4; b = a +1; p b';
Lexical::Analysis::scanner($a);
print "prog = $a\ntokens = @PL::Tutu::tokens\n";
```

Observa como la opción `-I..` hace que se busque por las librerías en el directorio padre del actual. Cuando ejecutamos `pruebalex.pl` obtenemos la lista de terminales:

```
$ ./pruebalex.pl
prog = int a,b; string c; c = "hello"; a = 4; b = a +1; p b
tokens = INT INT ID a PUN , ID b PUN ; STRING STRING ID c PUN ;
ID c PUN = STR hello PUN ; ID a PUN = NUM 4 PUN ; ID b PUN =
ID a PUN + NUM 1 PUN ; P P ID b
```

La última línea ha sido partida por razones de legibilidad, pero constituye una sola línea. Editemos el fichero `test.pl` en el directorio del módulo. Sus contenidos son como sigue:

```
$ cat -n test.pl
1  # Before 'make install' is performed this script should be runnable with
2  # 'make test'. After 'make install' it should work as 'perl test.pl'
3
4  #####
5
6  # change 'tests => 1' to 'tests => last_test_to_print';
7
8  use Test;
9  BEGIN { plan tests => 1 };
10 use PL::Tutu;
11 ok(1); # If we made it this far, we're ok.
12
13 #####
14
15 # Insert your test code below, the Test module is use()ed here so read
16 # its man page ( perldoc Test ) for help writing this test script.
17
```

En la línea 9 se establece el número de pruebas a realizar. La primera prueba aparece en la línea 11. Puede parecer que no es una prueba, ¡pero lo es!. Si se ha alcanzado la línea 11 es que se pudo cargar el módulo `PL::Tutu` y eso ¡tiene algún mérito!.

Seguiremos el consejo de la línea 15 y escribiremos nuestra segunda prueba al final del fichero `test.pl`:

```
$ cat -n test.pl | tail -7
16 # its man page ( perldoc Test ) for help writing this test script.
17
18 my $a = 'int a,b; string c; c = "hello"; a = 4; b = a +1; p b';
19 local @PL::Tutu::tokens = ();
20 Lexical::Analysis::scanner($a);
21 ok("@PL::Tutu::tokens" eq
22 'INT INT ID a PUN , ID b PUN ; STRING STRING ID c PUN ; ID c
   PUN = STR hello PUN ; ID a PUN = NUM 4 PUN ; ID b PUN =
   ID a PUN + NUM 1 PUN ; P P ID b');

```

La línea 22 ha sido partida por razones de legibilidad, pero constituye una sola línea. Ahora podemos ejecutar `make test` y comprobar que las dos pruebas funcionan:

```
$ make test
PERL_DL_NONLAZY=1 /usr/bin/perl -Ilib/arch -Ilib/lib -I/usr/lib/perl/5.6.1 \
-I/usr/share/perl/5.6.1 test.pl
1..2
ok 1
ok 2

```

¿Recordaste cambiar la línea 9 de `test.pl`? ¿Añadiste los nuevos ficheros a la lista en `MANIFEST`?

4.4.2. Práctica: Pruebas en el Análisis Léxico

Extienda su compilador para modularizar el analizador léxico tal y como se explicó en la sección 4.4.

1. Lea los siguientes documentos

- Test::Tutorial (Test::Tutorial - A tutorial about writing really basic tests) por Michael Schwern.
- Test Now, test Forever ” del libro de Scott [?].
- Perl Testing Reference Card por Ian Langworth.
- Chapter 4: Distributing Your Tests (and Code) del libro Perl Testing: A Developer's Notebook

2. Incluya la estrategia de pruebas de no regresión explicada en las secciones previas. Dado que ahora la estructura del terminal es una estructura de datos mas compleja (`token`, `[value, line_number]`) no podrá usar `is`, ya que este último sólo comprueba la igualdad entre escalares. Use `is_deeply` para comprobar que la estructura de datos devuelta por el analizador léxico es igual a la esperada. Sigue un ejemplo:

```
nereida:~/src/perl/YappWithDefaultAction/t> cat -n 15treeregswith2arrays.t
1  #!/usr/bin/perl -w
2  use strict;
3  #use Test::More qw(no_plan);
4  use Test::More tests => 3;
5  use_ok qw(Parse::Eyapp) or exit;

.. ..... etc., etc.

84 my $expected_tree = bless( {
85     'children' => [
86         bless( { 'children' => [
87             bless( { 'children' => [], 'attr' => 'a', 'token' => 'a' }, 'TERMINAL' )
88         ],
89         }, 'A' ),
90         bless( { 'children' => [
91             bless( { 'children' => [], 'attr' => 'c', 'token' => 'c' }, 'TERMINAL' )
92         ],
93         }, 'C' )
94     ],
95     }, 'ABC' );
96 is_deeply($t, $expected_tree, "deleting node between arrays");
```

3. Extienda los tests con una prueba en la que la entrada contenga un carácter ilegal. Obsérve que, tal y como esta escrito la rutina `scanner`, si la entrada tiene un carácter ilegal se ejecutarán las líneas

```
26     elsif (/\\G\\s*(.)/gc) {
27         Error::fatal "Caracter invalido: $1\\n";
28     }
```

lo que causa la parada del programa de prueba, al ejecutarse `fatal` el cuál llama a `croak`.

```
sub fatal {
    my $msg = join " ", @_ ;
    croak("Error: $msg\\n");
}
```

El objetivo es lograr que el programa de pruebas continúe ejecutando las subsiguientes pruebas. Para ello puede usar `Test::Exception` o bien `eval` y la variable especial `$@` para controlar que el programa `.t` no termine prematuramente. Repase las secciones sobre pruebas en [?].

4. Pruebe a dar como entrada un fichero vacío
5. Pruebe a dar como entrada un fichero que no existe
6. Pruebe a dar como entrada un fichero binario
7. Si tiene sentido en su caso, llame a las subrutinas con mas argumentos (y también con menos) de los que esperan.
8. Si tiene sentido en su caso, llame a las subrutinas con argumentos cuyo tipo no es el que se espera.

No use prototipos para lograrlo. No es una buena idea. Los prototipos en Perl a menudo producen un preprocesado del parámetro. Escriba código que controle que la naturaleza del parámetro es la que se espera. Por ejemplo:

```
sub tutu {
    my $refhash = shift;
    croak "Error" unless UNIVERSAL::isa($refhash, 'HASH');
    ...
}
```

9. Cambie los argumentos de orden (si es que se aplica a su código)
10. Comentarios: pruebe con `/* * */ a = 4; /* * / */`. Tambien con comentarios anidados (debería producirse un error)
11. Flotantes: compruebe su expresión regular con `0.0 0e0 .0 0 1.e-5 1.0e2 -2.0 .` (un punto sólo)
12. Cadenas. Pruebe con las cadenas

```
""
"h\"a\"h"
"\\"
```

Pruebe también con una cadena con varias líneas y otra que contenga un carácter de control en su interior.

13. Convierta los fallos (bugs) que encontró durante el desarrollo en pruebas
14. Compruebe la documentación usando el módulo `Test::Pod` de Andy Lester. Instálelo si es necesario.
15. Utilice el módulo `Test::Warn` para comprobar que los mensajes de warning (uso de **warn** and **carp**) se muestran correctamente.
16. Una prueba `SKIP` declara un bloque de pruebas que - bajo ciertas circunstancias - puede saltarse. Puede ser que sepamos que ciertas pruebas sólo funcionan en ciertos sistemas operativos o que la prueba requiera que ciertos paquetes están instalados o que la máquina disponga de ciertos recursos (por ejemplo, acceso a internet). En tal caso queremos que los tests se consideren si se dan las circunstancias favorables pero que en otro caso se descarten sin protestas. Consulte la documentación de los módulos `Test::More` y `Test::Harness` sobre pruebas tipo `SKIP`. El ejemplo que sigue declara un bloque de pruebas que pueden saltarse. La llamada a `skip` indica cuantos tests hay, bajo que condición saltarselos.

```

1  SKIP: {
2      eval { require HTML::Lint };
3
4      skip "HTML::Lint not installed", 2 if $@;
5
6      my $lint = new HTML::Lint;
7      isa_ok( $lint, "HTML::Lint" );
8
9      $lint->parse( $html );
10     is( $lint->errors, 0, "No errors found in HTML" );
11 }

```

Si el usuario no dispone del módulo `HTML::Lint` el bloque no será ejecutado. El módulo `Test::More` producirá oks que serán interpretados por `Test::Harness` como tests *skipped* pero *ok*.

Otra razón para usar una prueba `SKIP` es disponer de la posibilidad de saltarse ciertos grupos de pruebas. Por ejemplo, aquellas que llevan demasiado tiempo de ejecución y no son tan significativas que no se pueda prescindir de ellas cuando se introducen pequeños cambios en el código. El siguiente código muestra como usando una variable de entorno `TEST_FAST` podemos controlar que pruebas se ejecutan.

```

nereida:~/src/perl/YappWithDefaultAction/t> cat 02Cparser.t | head -n 56 -
#!/usr/bin/perl -w
use strict;
#use Test::More qw(no_plan);
use Test::More tests => 6;

use_ok qw(Parse::Eyapp) or exit;

SKIP: {
    skip "You decided to skip C grammar test (env var TEST_FAST)", 5 if $ENV{TEST_FAST} ;
    my ($grammar, $parser);
    $grammar=join(' ',<DATA>);
    $parser=new Parse::Eyapp(input => $grammar, inputfile => 'DATA', firstline => 52);

    #is($@, undef, "Grammar module created");

    # Does not work. May I have done s.t. wrong?
    #is(keys(%{$parser->{GRAMMAR}{NULLABLE}}), 43, "43 nullable productions");

    is(keys(%{$parser->{GRAMMAR}{NTERM}}), 233, "233 syntactic variables");

    is(scalar(@{$parser->{GRAMMAR}{UUTERM}}), 3, "3 UUTERM");

    is(scalar(keys(%{$parser->{GRAMMAR}{TERM}})), 108, "108 terminals");

    is(scalar(@{$parser->{GRAMMAR}{RULES}}), 825, "825 rules");

    is(scalar(@{$parser->{STATES}}), 1611, "1611 states");
}

__DATA__
/*
    This grammar is a stripped form of the original C++ grammar

```

from the GNU CC compiler :

YACC parser for C++ syntax.

Copyright (C) 1988, 89, 93-98, 1999 Free Software Foundation, Inc.

Hacked by Michael Tiemann (tiemann@cygnus.com)

The full gcc compiler and the original grammar file are freely available under the GPL license at :

<ftp://ftp.gnu.org/gnu/gcc/>

..... etc. etc.

*/

```
nereida:~/src/perl/YappWithDefaultAction> echo $TEST_FAST
```

```
1
```

```
nereida:~/src/perl/YappWithDefaultAction> make test
```

```
PERL_DL_NONLAZY=1 /usr/bin/perl "-MExtUtils::Command::MM" "-e" "test_harness(0, 'blib/lib
```

```
t/01calc.....ok
```

```
t/02Cparser.....ok
```

```
5/6 skipped: various reasons
```

```
t/03newgrammar.....ok
```

```
t/04foldandzero.....ok
```

```
t/05treewithvars.....ok
```

```
t/06meta.....ok
```

```
t/07translationschemetype.....ok
```

```
t/08tschemetypestar.....ok
```

```
t/09ts_with_defaultaction.....ok
```

```
t/10ts_with_treereg.....ok
```

```
etc., etc.....ok
```

```
t/28unshifttwoitems.....ok
```

```
t/29foldinglistsofexpressions.....ok
```

```
t/30complextreereg.....ok
```

```
t/32deletenodewithwarn.....ok
```

```
t/33moveinvariantoutofloop.....ok
```

```
t/34moveinvariantoutofloopcomplexformula....ok
```

```
All tests successful, 5 subtests skipped.
```

```
Files=33, Tests=113, 5 wallclock secs ( 4.52 cusr + 0.30 csys = 4.82 CPU)
```

Introduzca una prueba SKIP similar a la anterior y otra que si el módulo `Test::Pod` esta instalado comprueba que la documentación esta bien escrita. Estudie la documentación del módulo `Test::Pod`.

17. Introduzca pruebas TODO (que, por tanto, deben fallar) para las funciones que están por escribir (parser, Optimize, code_generator, transform). Repáse [?]. Sigue un ejemplo:

```
42 TODO: {
```

```
43     local $TODO = "Randomly generated problem";
```

```
44     can_ok('Algorithm::Knap01DP', 'GenKnap'); # sub GenKnap no ha sido escrita aún
```

```
45 }
```

18. Cuando compruebe el funcionamiento de su módulo *nunca descarte que el error pueda estar en el código de la prueba*. En palabras de Schwern

Code has bugs. Tests are code. Ergo, tests have bugs.

Michael Schwern

19. Instale el módulo Devel::Cover. El módulo Devel::Cover ha sido escrito por Paul Johnson y proporciona estadísticas del cubrimiento alcanzado por una ejecución. Para usarlo siga estos pasos:

```
pl@nereida:~/src/perl/YappWithDefaultAction$ cover -delete
Deleting database /home/pl/src/perl/YappWithDefaultAction/cover_db
pl@nereida:~/src/perl/YappWithDefaultAction$ HARNESS_PERL_SWITCHES=-MDevel::Cover make test
PERL_DL_NONLAZY=1 /usr/bin/perl "-MExtUtils::Command::MM" "-e" "test_harness(0, 'blib/lib
t/01calc.....ok
t/01calc.....ok
t/02Cparser.....ok
    5/6 skipped: various reasons
t/03newgrammar.....ok
t/03newgrammar.....ok
t/04foldandzero.....ok
etc., etc. ....ok
t/34moveinvariantoutofloopcomplexformula....ok
All tests successful, 5 subtests skipped.
Files=33, Tests=113, 181 wallclock secs (177.95 cusr + 2.94 csys = 180.89 CPU)
```

La ejecución toma ahora mucho mas tiempo: ¡181 segundos frente a los 5 que toma la ejecución sin cover !. Al ejecutar cover de nuevo obtenemos una tabla con las estadísticas de cubrimiento:

```
pl@nereida:~/src/perl/YappWithDefaultAction$ cover
Reading database from /home/pl/src/perl/YappWithDefaultAction/cover_db
```

File	stmt	bran	cond	sub	pod	time	total
blib/lib/Parse/Eyapp.pm	100.0	n/a	n/a	100.0	n/a	0.2	100.0
...lib/Parse/Eyapp/Driver.pm	72.4	63.2	50.0	64.3	0.0	21.3	64.4
...ib/Parse/Eyapp/Grammar.pm	90.9	77.8	66.7	100.0	0.0	16.6	84.3
blib/lib/Parse/Eyapp/Lalr.pm	91.4	72.6	78.6	100.0	0.0	48.3	85.6
blib/lib/Parse/Eyapp/Node.pm	74.4	58.3	29.2	88.2	0.0	1.6	64.7
...ib/Parse/Eyapp/Options.pm	86.4	50.0	n/a	100.0	0.0	2.7	72.8
...lib/Parse/Eyapp/Output.pm	82.3	47.4	60.0	70.6	0.0	3.7	70.0
.../lib/Parse/Eyapp/Parse.pm	100.0	n/a	n/a	100.0	n/a	0.2	100.0
...Parse/Eyapp/Treeregexp.pm	100.0	n/a	n/a	100.0	n/a	0.1	100.0
blib/lib/Parse/Eyapp/YATW.pm	89.4	63.9	66.7	85.7	0.0	4.8	77.6
...app/_TreeregexpSupport.pm	73.1	33.3	50.0	100.0	0.0	0.4	60.8
main.pm	52.2	0.0	n/a	80.0	0.0	0.0	45.7
Total	83.8	64.7	60.0	84.5	0.0	100.0	75.5

```
Writing HTML output to /home/pl/src/perl/YappWithDefaultAction/cover_db/coverage.html ...
pl@nereida:~/src/perl/YappWithDefaultAction$
```

El HTML generado nos permite tener una visión mas detallada de los niveles de cubrimiento.

Para mejorar el cubrimiento de tu código comienza por el informe de cubrimiento de subrutinas. Cualquier subrutina marcada como no probada es un candidato a contener errores o incluso a ser *código muerto*.

Para poder hacer el cubrimiento del código usando Devel::Cover, si se usa una `cs`h o `tc`sh se debe escribir:

```
nereida:~/src/perl/YappWithDefaultAction> setenv HARNESS_PERL_SWITCHES -MDevel::Cover
nereida:~/src/perl/YappWithDefaultAction> make test
PERL_DL_NONLAZY=1 /usr/bin/perl "-MExtUtils::Command::MM" "-e" "test_harness(0, 'blib/lib
t/01calc.....ok
t/01calc.....ok
t/02Cparser.....ok
    5/6 skipped: various reasons
t/03newgrammar.....ok
t/03newgrammar.....ok
t/04foldandzero.....ok
t/05treewithvars.....ok
t/06meta.....ok
t/06meta.....ok
t/07translationschemetype.....ok
.....ok
t/38tspostfix_resultisarray.....ok
t/39tspostfix.....ok
All tests successful, 5 subtests skipped.
Files=38, Tests=135, 210 wallclock secs (206.28 cusr + 3.27 csys = 209.55 CPU)
nereida:~/src/perl/YappWithDefaultAction>
```

Aún mas robusto - más independiente de la shell que usemos - es pasar las opciones en `HARNESS_PERL_SWITCHES` como parámetro a `make`:

```
make HARNESS_PERL_SWITCHES=-MDevel::Cover test
```

Añade el informe de cubrimiento al `MANIFEST` para que se incluya en la distribución que subas. Si lo consideras conveniente añade un directorio informes en los que vayan los informes asociados a esta práctica. Incluye en el `README` o en la documentación una breve descripción de donde están los informes.

20. Se conoce con el nombre de *perfilado* o *profiling* de un programa al estudio de su rendimiento mediante un programa (conocido como *profiler*) que monitoriza la ejecución del mismo mediante una técnica que interrumpe cada cierto tiempo el programa para comprobar en que punto de la ejecución se encuentra. Las estadísticas acumuladas se vuelcan al final de la ejecución en un fichero que puede ser visualizado mediante la aplicación apropiada.

En Perl hay dos módulos que permiten realizar profiling. El mas antiguo es `Devel::DProf`. La aplicación para visualizar los resultados se llama `dprofpp`. Sigue un ejemplo de uso:

```
nereida:~/src/perl/YappWithDefaultAction/t> perl -d:DProf 02Cparser.t
1..6
ok 1 - use Parse::Eyapp;
ok 2 - 233 syntactic variables
ok 3 - 3 UUTERM
ok 4 - 108 terminals
ok 5 - 825 rules
ok 6 - 1611 states
nereida:~/src/perl/YappWithDefaultAction/t> dprofpp tmon.out
Total Elapsed Time = 3.028396 Seconds
User+System Time = 3.008396 Seconds
```

Exclusive Times

%Time	ExclSec	CumulS	#Calls	sec/call	Csec/c	Name
31.4	0.945	1.473	1611	0.0006	0.0009	Parse::Eyapp::Lalr::_Transitions
17.5	0.528	0.528	1611	0.0003	0.0003	Parse::Eyapp::Lalr::_Closures
16.1	0.486	0.892	1	0.4861	0.8918	Parse::Eyapp::Lalr::_ComputeFollows
8.04	0.242	0.391	1	0.2419	0.3906	Parse::Yapp::Driver::_Parse
8.04	0.242	0.242	11111	0.0000	0.0000	Parse::Eyapp::Lalr::_ANON__
4.59	0.138	0.138	8104	0.0000	0.0000	Parse::Eyapp::Lalr::_Preds
2.66	0.080	0.080	1	0.0800	0.0800	Parse::Eyapp::Lalr::_SetDefaults
2.66	0.080	0.972	1	0.0800	0.9718	Parse::Eyapp::Lalr::_ComputeLA
2.46	0.074	0.074	3741	0.0000	0.0000	Parse::Eyapp::Parse::_Lexer
1.89	0.057	0.074	8310	0.0000	0.0000	Parse::Eyapp::Parse::_ANON__
0.96	0.029	0.028	1	0.0288	0.0276	Parse::Eyapp::Lalr::_SolveConflict
						s
0.66	0.020	0.050	6	0.0033	0.0083	Parse::Eyapp::Output::BEGIN
0.60	0.018	1.500	1	0.0176	1.4997	Parse::Eyapp::Lalr::_LRO
0.53	0.016	0.259	3	0.0054	0.0863	Parse::Eyapp::Lalr::_Digraph
0.33	0.010	0.010	1	0.0100	0.0100	Parse::Eyapp::Grammar::_SetNullable

Tambien es posible usar el módulo -MDevel::Profiler :

```
nereida:~/src/perl/YappWithDefaultAction/examples> perl -MDevel::Profiler eyapp 02Cparser
Unused terminals:
```

```
END_OF_LINE, declared line 128
ALL, declared line 119
PRE_PARSED_CLASS_DECL, declared line 120
```

27 shift/reduce conflicts and 22 reduce/reduce conflicts

```
nereida:~/src/perl/YappWithDefaultAction/examples> dprofpp tmon.out
```

Total Elapsed Time = 3.914144 Seconds

User+System Time = 3.917144 Seconds

Exclusive Times

%Time	ExclSec	CumulS	#Calls	sec/call	Csec/c	Name
22.3	0.877	1.577	1611	0.0005	0.0010	Parse::Eyapp::Lalr::_Transitions
17.8	0.700	0.700	1611	0.0004	0.0004	Parse::Eyapp::Lalr::_Closures
15.6	0.614	1.185	1	0.6142	1.1854	Parse::Eyapp::Lalr::_ComputeFollow
						s
9.60	0.376	0.545	1	0.3758	0.5453	Parse::Yapp::Driver::_Parse
7.99	0.313	0.313	8104	0.0000	0.0000	Parse::Eyapp::Lalr::_Preds
5.85	0.229	0.229	3	0.0763	0.0763	Parse::Eyapp::Lalr::_Digraph
4.06	0.159	0.159	3741	0.0000	0.0000	Parse::Eyapp::Parse::_Lexer
3.32	0.130	0.130	1	0.1300	0.1300	Parse::Eyapp::Lalr::DfaTable
2.27	0.089	0.089	1	0.0890	0.0890	Parse::Eyapp::Lalr::_SetDefaults
2.04	0.080	1.265	1	0.0800	1.2654	Parse::Eyapp::Lalr::_ComputeLA
1.17	0.046	0.057	1	0.0464	0.0567	Parse::Eyapp::Grammar::Rules
1.02	0.040	1.617	1	0.0397	1.6169	Parse::Eyapp::Lalr::_LRO
0.77	0.030	0.030	1185	0.0000	0.0000	Parse::Eyapp::Lalr::_FirstSfx
0.71	0.028	0.039	1	0.0284	0.0387	Parse::Eyapp::Grammar::RulesTable
0.54	0.021	0.021	1650	0.0000	0.0000	Parse::Eyapp::Grammar::classname

Presente un informe del perfil de su compilador. Añade el informe del perfil al MANIFEST para que se incluya en la distribución que subas.

21. El módulo `Devel::Size` proporciona la posibilidad de conocer cuanto ocupa una estructura de datos. Considere el siguiente ejemplo:

```
71 ..... codigo omitido
72
73 use Devel::Size qw(size total_size);
74 use Perl6::Form;
75
76 sub sizes {
77     my $d = shift;
78     my ($ps, $ts) = (size($d), total_size($d));
79     my $ds = $ts-$ps;
80     return ($ps, $ds, $ts);
81 }
82
83 print form(
84 ' =====',
85 '| VARIABLE | SOLO ESTRUCTURA |      SOLO DATOS |          TOTAL |',
86 '|-----+-----+-----+-----|',
87 '| $parser  | {>>>>>} bytes  | {>>>>>} bytes | {>>>>>} bytes |', sizes($parser),
88 '| $t       | {>>>>>} bytes  | {>>>>>} bytes | {>>>>>} bytes |', sizes($t),
89 ' =====',
90 );
```

Al ejecutarlo se obtiene esta salida:

..... salida previa omitida

```
=====
| VARIABLE | SOLO ESTRUCTURA |      SOLO DATOS |          TOTAL |
|-----+-----+-----+-----|
| $parser  |      748 bytes  |      991 bytes |     1739 bytes |
| $t       |       60 bytes  |     1237 bytes |     1297 bytes |
=====
```

Elabore un informe con el consumo de memoria de las variables mas importantes de su programa. Añadelo el informe al **MANIFEST** para que se incluya en la distribución que subas. Explica en el **README** o en la documentación el significado de los ficheros de informe.

4.4.3. Repaso: Pruebas en el Análisis Léxico

1. ¿Cuál es la diferencia entre los operadores `==` y `eq`?
2. ¿Cuáles son los parámetros de la función `ok`?
3. ¿Cuáles son los parámetros de la función `is`?
4. ¿Porqué es conveniente nombrar las pruebas con un nombre que empiece por un número?
5. ¿Como puedo ejecutar los tests en modo *verbose*?
6. ¿Como puedo probar un código que produce la detención del programa?
7. ¿Que contiene la variable `$@`?
8. ¿Que hace la función `like`?

9. ¿Que contiene la variable `$_?`? ¿Y `$+?` (Consulte 3.1.4)
10. ¿Porqué la función `use_ok` es llamada dentro de un `BEGIN`?
11. ¿Que es una prueba `SKIP`?
12. ¿Que es una prueba `TODO`?
13. ¿Que hace la función `pod_file_ok`? ¿A que módulo pertenece?
14. ¿Que hace el operador `tr`?
15. ¿Qué devuelve el operador `tr`?
16. ¿Que hace la opción `d` de `tr`? (consulte 3.5)
17. Explique la conducta de la siguiente sesión con el depurador:

4.5. Conceptos Básicos para el Análisis Sintáctico

Suponemos que el lector de esta sección ha realizado con éxito un curso en teoría de autómatas y lenguajes formales. Las siguientes definiciones repasan los conceptos mas importantes.

Definición 4.5.1. Dado un conjunto A , se define A^* el cierre de Kleene de A como: $A^* = \bigcup_{n=0}^{\infty} A^n$

Se admite que $A^0 = \{\epsilon\}$, donde ϵ denota la palabra vacía, esto es la palabra que tiene longitud cero, formada por cero símbolos del conjunto base A .

Definición 4.5.2. Una gramática G es una cuaterna $G = (\Sigma, V, P, S)$. Σ es el conjunto de terminales. V es un conjunto (disjunto de Σ) que se denomina conjunto de variables sintácticas o categorías gramaticales, P es un conjunto de pares de $V \times (V \cup \Sigma)^*$. En vez de escribir un par usando la notación $(A, \alpha) \in P$ se escribe $A \rightarrow \alpha$. Un elemento de P se denomina producción. Por último, S es un símbolo del conjunto V que se denomina símbolo de arranque.

Definición 4.5.3. Dada una gramática $G = (\Sigma, V, P, S)$ y $\mu = \alpha A \beta \in (V \cup \Sigma)^*$ una frase formada por variables y terminales y $A \rightarrow \gamma$ una producción de P , decimos que μ deriva en un paso en $\alpha \gamma \beta$. Esto es, derivar una cadena $\alpha A \beta$ es sustituir una variable sintáctica A de V por la parte derecha γ de una de sus reglas de producción. Se dice que μ deriva en n pasos en δ si deriva en $n - 1$ pasos en una cadena $\alpha A \beta$ la cual deriva en un paso en δ . Se escribe entonces que $\mu \xRightarrow{*} \delta$. Una cadena deriva en 0 pasos en si misma.

Definición 4.5.4. Dada una gramática $G = (\Sigma, V, P, S)$ se denota por $L(G)$ o lenguaje generado por G al lenguaje:

$$L(G) = \{x \in \Sigma^* : S \xRightarrow{*} x\}$$

Esto es, el lenguaje generado por la gramática G esta formado por las cadenas de terminales que pueden ser derivados desde el símbolo de arranque.

Definición 4.5.5. Una derivación que comienza en el símbolo de arranque y termina en una secuencia formada por sólo terminales de Σ se dice completa.

Una derivación $\mu \xRightarrow{*} \delta$ en la cual en cada paso $\alpha A x$ la regla de producción aplicada $A \rightarrow \gamma$ se aplica en la variable sintáctica mas a la derecha se dice una derivación a derechas

Una derivación $\mu \xRightarrow{*} \delta$ en la cual en cada paso $x A \alpha$ la regla de producción aplicada $A \rightarrow \gamma$ se aplica en la variable sintáctica mas a la izquierda se dice una derivación a izquierdas

Definición 4.5.6. Observe que una derivación puede ser representada como un árbol cuyos nodos están etiquetados en $V \cup \Sigma$. La aplicación de la regla de producción $A \rightarrow \gamma$ se traduce en asignar como hijos del nodo etiquetado con A a los nodos etiquetados con los símbolos $X_1 \dots X_n$ que constituyen la frase $\gamma = X_1 \dots X_n$. Este árbol se llama árbol sintáctico concreto asociado con la derivación.

Definición 4.5.7. Observe que, dada una frase $x \in L(G)$ una derivación desde el símbolo de arranque da lugar a un árbol. Ese árbol tiene como raíz el símbolo de arranque y como hojas los terminales $x_1 \dots x_n$ que forman x . Dicho árbol se denomina árbol de análisis sintáctico concreto de x . Una derivación determina una forma de recorrido del árbol de análisis sintáctico concreto.

Definición 4.5.8. Una gramática G se dice ambigua si existe alguna frase $x \in L(G)$ con al menos dos árboles sintácticos. Es claro que esta definición es equivalente a afirmar que existe alguna frase $x \in L(G)$ para la cual existen dos derivaciones a izquierda (derecha) distintas.

4.5.1. Ejercicio

Dada la gramática con producciones:

```

program → declarations statements | statements
declarations → declaration ';' declarations | declaration ';'
declaration → INT idlist | STRING idlist
statements → statement ';' statements | statement
statement → ID '=' expression | P expression
expression → term '+' expression | term
term → factor '*' term | factor
factor → '(' expression ')' | ID | NUM | STR
idlist → ID ',' idlist | ID

```

En esta gramática, Σ está formado por los caracteres entre comillas simples y los símbolos cuyos identificadores están en mayúsculas. Los restantes identificadores corresponden a elementos de V . El símbolo de arranque es $S = \text{program}$.

Conteste a las siguientes cuestiones:

1. Describa con palabras el lenguaje generado.
2. Construya el árbol de análisis sintáctico concreto para cuatro frases del lenguaje.
3. Señale a que recorridos del árbol corresponden las respectivas derivaciones a izquierda y a derecha en el apartado 2.
4. ¿Es ambigua esta gramática?. Justifique su respuesta.

4.6. Análisis Sintáctico Predictivo Recursivo

La siguiente fase en la construcción del analizador es la fase de análisis sintáctico. Esta toma como entrada el flujo de terminales y construye como salida el árbol de análisis sintáctico abstracto.

El árbol de análisis sintáctico abstracto es una representación compactada del árbol de análisis sintáctico concreto que contiene la misma información que éste.

Existen diferentes métodos de análisis sintáctico. La mayoría caen en una de dos categorías: ascendentes y descendentes. Los ascendentes construyen el árbol desde las hojas hacia la raíz. Los descendentes lo hacen en modo inverso. El que describiremos aquí es uno de los mas sencillos: se denomina método de análisis predictivo descendente recursivo.

4.6.1. Introducción

En este método se asocia una subrutina con cada variable sintáctica $A \in V$. Dicha subrutina (que llamaremos **A**) reconocerá el lenguaje generado desde la variable A :

$$L_A(G) = \{x \in \Sigma^* : A \xRightarrow{*} x\}$$

En este método se escribe una rutina **A** por variable sintáctica $A \in V$. Se le da a la rutina asociada el mismo nombre que a la variable sintáctica asociada. La función de la rutina **A** asociada con la variable $A \in V$ es reconocer el lenguaje $L(A)$ generado por A . La estrategia general que sigue la rutina **A** para reconocer $L(A)$ es decidir en términos del terminal a en la entrada que regla de producción concreta $A \rightarrow \alpha$ se aplica para a continuación comprobar que la entrada que sigue pertenece al lenguaje generado por α . En un analizador predictivo descendente recursivo (APDR) se asume que el símbolo que actualmente está siendo observado (denotado **lookahead**) permite determinar unívocamente que producción de A hay que aplicar. Una vez que se ha determinado que la regla por la que continuar la derivación es $A \rightarrow \alpha$ se procede a reconocer $L_\alpha(G)$, el lenguaje generado por α . Si $\alpha = X_1 \dots X_n$, las apariciones de terminales X_i en α son emparejadas con los terminales en la entrada mientras que las apariciones de variables $X_i = B$ en α se traducen en llamadas a la correspondiente subrutina asociada con **B**.

Para ilustrar el método, simplificaremos la gramática presentada en el ejercicio 4.5.1 eliminando las declaraciones:

```

statements → statement ';' statements | statement
statement → ID '=' expression | P expression
expression → term '+' expression | term
term → factor '*' term | factor
factor → '(' expression ')' | ID | NUM

```

La secuencia de llamadas cuando se procesa la entrada mediante el siguiente programa construye “implícitamente” el árbol de análisis sintáctico concreto.

Dado que estamos usando **strict** se requiere prototipar las funciones al comienzo del fichero:

```

sub parse();
sub statements();
sub statement();
sub expression();
sub term();
sub factor();
sub idlist();
sub declaration();
sub declarations();

```

Para saber mas sobre prototipos consulte [?].

Programa 4.6.1. 1 sub match {

```

2  my $t = shift;
3
4  if ($lookahead eq $t) {
5      ($lookahead, $value) = splice @tokens,0,2;
6      if (defined($lookahead)) {
7          $lookahead = $value if ($lookahead eq 'PUN');
8      } else { $lookahead = 'EOI'; }
9  }
10 else { error("Se esperaba $t y se encontro $lookahead\n"); }
11 }
12
13 sub statement {
14     if ($lookahead eq 'ID') { match('ID'); match('='); expression; }
15     elsif ($lookahead eq 'P') { match('P'); expression; }
16     else { error('Se esperaba un identificador'); }
17 }
18
19 sub term() {
20     factor;
21     if ($lookahead eq '*') { match('*'); term; }
22 }
23
24 sub expression() {
25     term;
26     if ($lookahead eq '+') { match('+'); expression; }
27 }
28
29 sub factor() {
30     if ($lookahead eq 'NUM') { match('NUM'); }
31     elsif ($lookahead eq 'ID') { match('ID'); }
32     elsif ($lookahead eq '(') { match('('); expression; match(')'); }
33     else { error("Se esperaba (, NUM o ID"); }

```



```

34 }
35
36 sub statements {
37     statement;
38     if ($lookahead eq ';' ) { match(';'); statements; }
39 }
40
41 sub parser {
42     ($lookahead, $value) = splice @tokens,0,2;
43     statements; match('EOI');
44 }

```

Como vemos en el ejemplo, el análisis predictivo confía en que, si estamos ejecutando la entrada del procedimiento A , el cuál está asociado con la variable $A \in V$, el símbolo terminal que esta en la entrada a determine de manera unívoca la regla de producción $A \rightarrow a\alpha$ que debe ser procesada.

Si se piensa, esta condición requiere que todas las partes derechas α de las reglas $A \rightarrow \alpha$ de A “comiencen” por diferentes símbolos. Para formalizar esta idea, introduciremos el concepto de conjunto $FIRST(\alpha)$:

Definición 4.6.1. Dada una gramática $G = (\Sigma, V, P, S)$ y un símbolo $\alpha \in (V \cup \Sigma)^*$ se define el conjunto $FIRST(\alpha)$ como:

$$FIRST(\alpha) = \left\{ b \in \Sigma : \alpha \xRightarrow{*} b\beta \right\} \cup N(\alpha)$$

donde:

$$N(\alpha) = \begin{cases} \{\epsilon\} & \text{si } \alpha \xRightarrow{*} \epsilon \\ \emptyset & \text{en otro caso} \end{cases}$$

Podemos reformular ahora nuestra afirmación anterior en estos términos: Si $A \rightarrow \gamma_1 \mid \dots \mid \gamma_n$ y los conjuntos $FIRST(\gamma_i)$ son disjuntos podemos construir el procedimiento para la variable A siguiendo este pseudocódigo:

```

sub A {
    if ($lookahead in FIRST(gamma_1)) { imitar gamma_1 }
    elsif ($lookahead in FIRST(gamma_2)) { imitar gamma_2 }
    ...
    else ($lookahead in FIRST(gamma_n)) { imitar gamma_n }
}

```

Donde si γ_j es $X_1 \dots X_k$ el código `gamma_j` consiste en una secuencia $i = 1 \dots k$ de llamadas de uno de estos dos tipos:

- Llamar a la subrutina `X_i` si X_i es una variable sintáctica
- Hacer una llamada a `match(X_i)` si X_i es un terminal

4.6.2. Ejercicio: Recorrido del árbol en un ADPR

¿En que forma es recorrido el árbol de análisis sintáctico concreto en un analizador descendente predictivo recursivo? ¿En que orden son visitados los nodos?

4.6.3. Ejercicio: Factores Comunes

En el programa 4.6.1 el reconocimiento de las categorías gramaticales `statements`, `expression` y `term` (líneas 19-27) difiere del resto. Observe las reglas:

```

statements → statement ';' statements | statement
expression → term '+' expression | term
term → factor '*' term | factor

```

¿Son disjuntos los conjuntos $FIRST(\gamma_i)$ para las partes derechas de las reglas de statements? ¿Son disjuntos los conjuntos $FIRST(\gamma_i)$ para las partes derechas de las reglas de expression? ¿Son disjuntos los conjuntos $FIRST(\gamma_i)$ para las partes derechas de las reglas de term?

Si se tiene una variable con producciones:

$$A \rightarrow \alpha\beta \mid \alpha\gamma$$

Las dos producciones tienen un *máximo factor común* en la izquierda de su parte derecha α . Asumimos que $FIRST(\beta) \cap FIRST(\gamma) = \emptyset$.

1. ¿Cómo puede modificarse la gramática para obtener una nueva gramática que cumpla la condición de que las partes derechas tienen conjuntos $FIRST(\gamma_i)$ disjuntos?
2. ¿Puede modificarse la técnica APDR para que funcione sobre gramáticas con este tipo de producciones?. Observe el código asociado con statements, expression y term. ¿Cómo sería el esquema general?

4.6.4. Derivaciones a vacío

Surge un problema cuando $A \rightarrow \gamma_1 \mid \dots \mid \gamma_n$ y la palabra vacía está en alguno de los conjuntos $FIRST(\gamma_i)$. ¿Que hacer entonces?

Nótese que si $A \rightarrow \gamma$ y $\epsilon \in FIRST(\gamma)$ es porque existe una derivación $\gamma \xRightarrow{*} \epsilon$. ¿Que terminales podemos legalmente encontrarnos cuando estamos en la subrutina **A**? Consideremos una derivación desde el símbolo de arranque en la que se use la producción $A \rightarrow \gamma$. Dicha derivación forzosamente tendrá la forma:

$$S \xRightarrow{*} \beta A a\mu \xRightarrow{*} \beta \gamma a\mu \xRightarrow{*} \beta a\mu.$$

Cualquier terminal $a \in \Sigma$ que pueda aparecer en una derivación desde el símbolo de arranque inmediatamente a continuación de la variable A es susceptible de ser visto cuando se esta analizando A y se aplicó $A \rightarrow \gamma$ con $\gamma \xRightarrow{*} \epsilon$. Esto nos lleva a la definición del conjunto $FOLLOW(A)$ como conjunto de terminales que pueden aparecer a continuación de A en una derivación desde el símbolo de arranque:

Definición 4.6.2. Dada una gramática $G = (\Sigma, V, P, S)$ y una variable $A \in V$ se define el conjunto $FOLLOW(A)$ como:

$$FOLLOW(A) = \left\{ b \in \Sigma : \exists S \xRightarrow{*} \alpha A b \beta \right\} \cup E(A)$$

donde

$$E(A) = \begin{cases} \{\$ \} & \text{si } S \xRightarrow{*} \alpha A \\ \emptyset & \text{en otro caso} \end{cases}$$

Aquí $\$$ denota el final de la entrada (que se corresponde en el código Perl anterior con el terminal EOI).

Si $A \rightarrow \gamma_1 \mid \dots \mid \gamma_n$ dado que los conjuntos $FIRST(\gamma_i)$ han de ser disjuntos para que un analizador predictivo APDR funcione, sólo una parte derecha puede contener la palabra vacía en su $FIRST$. Supongamos que es γ_n . Podemos reformular la construcción del procedimiento para la variable A siguiendo este pseudocódigo:

```
sub A {
  if ($lookahead in FIRST(gamma_1)) { imitar gamma_1 }
  elsif ($lookahead in FIRST(gamma_2)) { imitar gamma_2 }
  ...
  else ($lookahead in FIRST(gamma_n) or $lookahead in FOLLOW(A)) { imitar gamma_n }
}
```

Un caso particular de $\gamma_n \xRightarrow{*} \epsilon$ es que $\gamma_n = \epsilon$. En tal caso, y como es obvio, el significado de `imitar gamma_n` es equivalente a ejecutar una sentencia vacía.

4.6.5. Construcción de los conjuntos de Primeros y Siguientes

Algoritmo 4.6.1. *Construcción de los conjuntos $FIRST(X)$*

Repita el siguiente conjunto de reglas hasta que no se puedan añadir mas símbolos terminales o a ningún conjunto $FIRST(X)$:

1. Si $X \in \Sigma$ entonces $FIRST(X) = X$
2. Si $X \rightarrow \epsilon$ entonces $FIRST(X) = FIRST(X) \cup \{\epsilon\}$
3. Si $X \in V$ y $X \rightarrow Y_1 Y_2 \cdots Y_k \in P$ entonces

$i = 1;$
 hacer
 $FIRST(X) = FIRST(X) \cup FIRST^*(Y_i);$
 $i++;$
 mientras ($i \leq k$ y $\epsilon \in FIRST(Y_i)$)

4. Añadir ϵ a $FIRST(X)$ si $i \geq k$ y $\epsilon \in FIRST(Y_k)$

Aquí $FIRST^*(Y)$ denota al conjunto $FIRST(Y) - \{\epsilon\}$.

Este algoritmo puede ser extendido para calcular $FIRST(\alpha)$ para $\alpha = X_1 X_2 \cdots X_n \in (V \cup \Sigma)^*$. El esquema es análogo al de un símbolo individual.

Algoritmo 4.6.2. *Construcción del conjunto $FIRST(\alpha)$*

Repita siguiente conjunto de reglas hasta que no se puedan añadir mas símbolos terminales o a ningún conjunto $FIRST(\alpha)$:

$i = 1;$
 $FIRST(\alpha) = \emptyset;$
 hacer
 $FIRST(\alpha) = FIRST(\alpha) \cup FIRST^*(X_i);$
 $i++;$
 mientras ($i \leq n$ y $\epsilon \in FIRST(X_i)$)

Algoritmo 4.6.3. *Construcción de los conjuntos $FOLLOW(A) \forall A \in V$:*

Repetir los siguientes pasos hasta que ninguno de los conjuntos $FOLLOW$ cambie:

1. $FOLLOW(S) = \{\$ \}$ ($\$$ representa el final de la entrada)
2. Si $A \rightarrow \alpha B \beta$ entonces

$$FOLLOW(B) = FOLLOW(B) \cup (FIRST(\beta) - \{\epsilon\})$$

3. Si $A \rightarrow \alpha B$ o $A \rightarrow \alpha B \beta$ y $\epsilon \in FIRST(\beta)$ entonces

$$FOLLOW(B) = FOLLOW(B) \cup FOLLOW(A)$$

4.6.6. Ejercicio: Construir los $FIRST$

Construya los conjuntos $FIRST$ de las partes derechas de las reglas de producción de la gramática presentada en el ejercicio 4.5.1.

4.6.7. Ejercicio: Calcular los *FOLLOW*

Modificamos la gramática de la sección 4.6.1 para que admita la sentencia vacía:

```
statements → statement ';' statements | statement
statement → ID '=' expression | P expression | ε
expression → term '+' expression | term
term → factor '*' term | factor
factor → '(' expression ')' | ID | NUM
```

Calcule los conjuntos *FOLLOW*. ¿Es la nueva gramática susceptible de ser analizada por un analizador predictivo descendente recursivo? ¿Cómo sería el código para la subrutina `statements`?. Escríbalo.

4.6.8. Práctica: Construcción de los *FIRST* y los *FOLLOW*

He escrito un módulo llamado Grammar que provee la función `Grammar::Parse` la cual recibe una cadena conteniendo la gramática en formato `yacc` o `eyapp` y devuelve una referencia a un hash conteniendo la información pertinente para el tratamiento de la gramática. Para instalar el módulo tenga en cuenta que depende del módulo `Parse::Yapp`.

Para ilustrar el uso vea los ejemplos en el directorio `scripts`. En concreto veamos el programa `grammar.pl`.

```
Grammar/scripts$ cat -n grammar.pl
1  #!/usr/bin/perl -w -I../lib
2  use strict;
3  use Grammar;
4  use Data::Dumper;
5
6  sub usage {
7      print <<"EOI";
8      usage:
9      $0 input_grammar
10  EOI
11      die "\n";
12  }
13
14  usage() unless @ARGV;
15  my $filename = shift;
16
17  local $/ = undef;
18  open my $FILE, "$filename";
19  my $grammar = <$FILE>;
20  my $x = Grammar::Parse($grammar);
21
22  print Dumper($x);
```

Vamos a darle como entrada la gramática en el fichero `aSb.yyp` conteniendo una gramática:

```
Grammar/scripts$ cat -n aSb.yyp
1  %%
2  S:
3      |   'a' S 'b'
4  ;
5  %%
```

Las gramáticas aceptadas por `Grammar::Parse` se adaptan a la sintaxis de las gramáticas reconocidas por `Parse::Yapp`. Una gramática (normalmente con tipo `.yp`) consta de tres partes: la cabeza, el cuerpo y la cola. Cada una de las partes va separada de las otras por el símbolo `%%` en una línea aparte. Así, el `%%` de la línea 1 separa la cabeza del cuerpo. En la cabecera se colocan las declaraciones de terminales (directiva `%token`), cual es el símbolo de arranque (directiva `%start`), etc. El cuerpo contiene las reglas de la gramática y las acciones asociadas. Por último, la cola en nuestro caso no es usada y es vacía. En general, la cola contiene las rutinas de soporte al código que aparece en las acciones así como, posiblemente, rutinas para el análisis léxico y el tratamiento de errores.

La salida de `Grammar::Parse` es una referencia a un hash cuyas entradas vienen explicadas por los comentarios.

```
Grammar/scripts$ grammar.pl aSb.yp
$VAR1 = {
  'SYMS' => { 'S' => 2, '"b"' => 3, '"a"' => 3 }, # Símbolo => línea
  'NULL' => { 'S' => 1 }, # símbolos que se anulan
  'RULES' => [
    [ 'S', [] ], # S produce vacío
    [ 'S', [ '"a"', 'S', '"b"' ] ] # S -> aSb
  ],
  'START' => 'S', # Símbolo de arranque
  'TERM' => [ '"b"', '"a"' ], # terminales /tokens
  'NTERM' => { 'S' => [ 0, 1 ] } # índices de las reglas de las variables sintácticas
};
```

Usando la estructura devuelta por la función `Grammar::Parse` escriba un módulo que provea funciones para computar los `FIRST` y los `FOLLOW` de las variables sintácticas de la gramática. No olvide escribir la documentación. Incluya una prueba por cada una de las gramáticas que figuran en el directorio `scripts` del módulo `Grammar`.

Puede encontrar la práctica *casi hecha* en `PL::FirstFollow`. Asegúrese de entender el algoritmo usado. Aumente el número de pruebas y haga un análisis de cubrimiento.

4.6.9. Gramáticas LL(1)

Una gramática $G = (\Sigma, V, P, S)$ cuyo lenguaje generado $L(G)$ puede ser analizado por un analizador sintáctico descendente recursivo predictivo se denomina *LL(1)*. Una gramática es LL(1) si y sólo si para cualesquiera dos producciones $A \rightarrow \alpha$ y $A \rightarrow \beta$ de G se cumple:

1. $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$
2. Si $\epsilon \in FIRST(\alpha)$, entonces $FIRST(\alpha) \cap FOLLOW(A) = \emptyset$

¿De donde viene el nombre LL(1)? La primera L hace alusión al hecho de que el flujo de terminales se lee de izquierda a derecha, accediendo a la entrada por su izquierda (*Left*). La segunda L se refiere a que el método de análisis predictivo construye una derivación a izquierdas. El número entre paréntesis indica el número de terminales que debemos consultar para decidir que regla de producción se aplica. Así, en una gramática LL(2) la decisión final de que producción elegir se hace consultando los dos terminales a la entrada.

4.6.10. Ejercicio: Caracterización de una gramática LL(1)

Cuando se dice que una gramática es LL(1) si, y sólo si:

1. $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$
2. Si $\epsilon \in FIRST(\alpha)$, entonces $FIRST(\alpha) \cap FOLLOW(A) = \emptyset$

se asume que los conjuntos $FIRST(\alpha)$ no son vacíos.

- ¿Que se puede decir de la regla $A \rightarrow \alpha$ si $FIRST(\alpha) = \emptyset$?
- ¿Que se puede decir de la variable A si $FOLLOW(A) = \emptyset$?

4.6.11. Ejercicio: Ambigüedad y LL(1)

¿Puede una gramática LL(1) ser ambigua?. Razone su respuesta.

4.6.12. Práctica: Un analizador APDR

Siguiendo con la construcción del compilador para el lenguaje Tutu, escriba un analizador APDR para la siguiente gramática. Reutilice el código de las prácticas de las secciones anteriores (4.3 y 4.4).

```

program → declarations statements | statements
declarations → declaration ';' declarations | declaration ';'
declaration → INT idlist | STRING idlist
statements → statement ';' statements | statement
statement → ID '=' expression | P expression | ε
expression → term '+' expression | term
term → factor '*' term | factor
factor → '(' expression ')' | ID | NUM | STR
idlist → ID ',' idlist | ID

```

4.6.13. Práctica: Generación Automática de Analizadores Predictivos

Objetivo Escriba un módulo `GAP.pm` que provea una subrutina `gap` para la generación automática de un APDR supuesto que la gramática de entrada es LL(1).

La subrutina `gap` recibe como entrada la gramática según la estructura de datos generada por la función `Grammar::Parse` de la versión 0.3 del módulo `Grammar`.

El Módulo Grammar La estructura de datos generada por la función `Grammar::Parse` se explicó en la práctica 4.6.8. La estructura ha sido extendida en esta versión para incluir el código que se sitúa en la zona de cola. Por ejemplo, dada la gramática de entrada:

```

Grammar/03/scripts$ cat -n aSb.y
 1  %%
 2  S:
 3      |   'a' S 'b'
 4      ;
 5  %%
 6
 7  sub Lex {
 8      local $_ = shift; # input
 9      my @tokens;
10
11
12      while ($_) {
13          s/^\s*//; # fuera blancos
14          push @tokens, $1, $1 if s/^(.)//s
15      }
16      @tokens;
17  }
18
19  sub main {

```

```

20   my $filename = shift;
21   my $input;
22
23   if (defined($filename)) {
24       local $/ = undef;
25       open my $FILE, $filename or die "No se pudo abrir $filename\n";
26       $input = <$FILE>;
27       close($FILE);
28   }
29   else { $input = <STDIN> }
30
31   my @tokens = Lex($input);
32   Parse(@tokens); # Llamada al analizador generado
33   print "Sintácticamente correcto\n";
34 }

```

se genera la siguiente estructura de datos:

```

{
'SYMS' => { 'S' => 2, 'b' => 3, 'a' => 3 }, # Símbolo => línea de aparición
'NULL' => { 'S' => 1 }, # Símbolos que se anulan
'RULES' => [ # Reglas
    [ 'S', [] ], # S produce vacío
    [ 'S', [ 'a', 'S', 'b' ] ] # S-> a S b
],
'START' => 'S', # Símbolo de arranque
'TERM' => [ 'b', 'a' ], # Terminales
'NTERM' => { 'S' => [ 0, 1 ] } # Variables sintácticas e índices de las reglas de esa variable
'TAIL' => [ # [ 'Código de cola', línea en la que está el segundo %% ]
,

sub Lex {
    local $_ = shift; # input
    my @tokens;

    while ($_) {
        s/^\s*//; # fuera blancos
        push @tokens, $1, $1 if s/^(.)//s
    }
    @tokens;
}

sub main {
    my $filename = shift;
    my $input;

    if (defined($filename)) {
        local $/ = undef;
        open my $FILE, $filename or die "No se pudo abrir $filename\n";
        $input = <$FILE>;
        close($FILE);
    }
    else { $input = <STDIN> }
}

```

```

my @tokens = Lex($input);
my $ok = Parse(@tokens); # Llamada al analizador generado
print "Sintácticamente correcto\\n" if $ok;
}

', 5 ], # línea en la que está el segundo %%
};

```

Así pues la entrada con clave `TAIL` contiene el código auxiliar de cola. Este código debe ser incluido por su programa dentro del texto del paquete generado por `gap`.

Descripción del objetivo: La función `gap` La función `gap` también recibe como entrada el nombre del package:

```
$package_text = &gap($grammar, 'Package_name');
```

La función `gap` retorna una cadena conteniendo el `package` en el que están las subrutinas del analizador sintáctico.

La idea es que dicha cadena se salvará en un fichero con nombre `Package_name.pm` que podrá posteriormente ser usado (`use Package_name`) por un programa que necesite analizar entradas que se conforman de acuerdo a la especificación de la gramática.

Descripción del objetivo: La función `parser` La rutina principal del paquete generado se ha de llamar `parser` (esto es, su nombre completo es: `Package_name::parser`. Evidentemente `Package_name` debe ser un nombre Perl válido). Ninguna subrutina deberá ser exportada sino que deberán ser llamadas por su nombre completo.

La subrutina `parser` recibe como argumento el array de terminales, obtiene el primer terminal y llama a la subrutina asociada con el símbolo de arranque. Los terminales están representados como parejas (*terminal, atributo*).

Observe que, una vez que la cadena `$package_text` conteniendo el paquete ha sido generada y salvada en un fichero con nombre `Package_name.pm`, podemos escribir un programa cliente:

```

use strict;
use Package_name;

&Package_name::main;

```

Este programa espera una entrada desde fichero o `STDIN` e informa si dicha entrada es sintácticamente correcta o no para la gramática en cuestión.

Cálculo de los `First` y los `Follow` con `PL::FirstFollow` Para facilitar la escritura de `GAP.pm` pueden hacer uso del módulo `PL::FirstFollow` el cual calcula los *FIRST* y los *FOLLOW*. El módulo `PL::FirstFollow` depende de `Set::Scalar` escrito por Jarkko Hietaniemi: instálelo primero.

Deberá familiarizarse con `PL::FirstFollow`, rellenar la documentación de todas las subrutinas (apariciones de `????` en el texto) y escribir la documentación siguiendo el template que se provee. *Rellene los fragmentos de código que se han sustituido por signos de interrogación.* Haga un estudio de cubrimiento y añada pruebas para mejorar el actual. El actual cubrimiento es:

File	stmt	bran	cond	sub	pod	time	total
...ammar-0.03/lib/Grammar.pm	100.0	n/a	n/a	100.0	0.0	75.3	97.2
blib/lib/PL/FirstFollow.pm	100.0	92.9	50.0	100.0	0.0	24.7	95.1
Total	100.0	92.9	50.0	100.0	0.0	100.0	95.5

Si observa un fallo en `PL::FirstFollow` háganoslo saber y además de resolverlo escriba una prueba para detectar el fallo.

Haga un estudio de profiling de su aplicación.

Uso de Templates Un módulo que puede facilitar la escritura de esta práctica es `Text::Template` debido a Mark Jason Dominus. El siguiente ejemplo de uso es un fragmento de un traductor - que nunca acabo de terminar - que toma con fuente un fichero en el formato que usa Moodle para los cuestionarios (conocido como formato GIFT) y lo convierte en un cuestionario `LATEX`:

```
lhp@nereida:~/projects/Gift2LaTeX/Gift2LaTeX/lib$ cat -n Gift2LaTeX.pm
 1  package Gift2LaTeX;
 2
 3  use strict;
 4  use warnings;
 5  use Gift;
 6  use Text::Template;
 7  use HTML::Latex;
..
49  package Gift::TRUEFALSE; # True-false questions belong to this class
50
51  { # closure
52
53      die "Can't find $TEMPLATE_DIR/TRUEFALSE_question.tep\n"
54          unless -e "$TEMPLATE_DIR/TRUEFALSE_question.tep";
55      my $tfq_tmpl = Text::Template->new( #tfq = true-false question
56          DELIMITERS => ['%<', '%>'];
57          SOURCE => "$TEMPLATE_DIR/TRUEFALSE_question.tep",
58      );
..
67  sub gen_latex {
68      my $self = shift;
69
70      ##### Generate latex for question
71      my $prefix = $self->PREFIX;
72
73      my $sufix = $self->POSTSTATE;
74
75      $self->Error("Only HTML and PLAIN formats are supported\n")
76          unless (! $self->FORMAT or ($self->FORMAT =~ m{html|plain}i));
77
78      my ($prefix_tex, $sufix_tex);
79      if (defined($self->FORMAT) and $self->FORMAT =~ m{plain}i) {
80          $prefix_tex = $prefix;
81          $sufix_tex = $sufix;
82      }
83      else { # HTML by default
..
86  }
87  my $params = {
88      prefix => $prefix_tex,
89      sufix  => $sufix_tex,
90      separator => $separator,
91      label   => $label_prefix.$question_number,
92      question_number => $question_number
```

```

93     };
94     my $question_tex = $tfq_tmpl->fill_in(HASH => $params);
95     ##### Generate latex for answer
96     .....
101 }
102 }

```

En la línea 55 se crea el template. El template se lee desde el fichero "\$TEMPLATE_DIR/TRUEFALSE_question.tep" cuyo contenido es una mezcla de texto (en este caso texto L^AT_EX y HTML) con código Perl: El código Perl aparece acotado entre los delimitadores '%<' y '%>'.

```

lhp@nereida:~/projects/Gift2LaTeX/Gift2LaTeX/etc/en$ cat -n TRUEFALSE_question.tep

```

```

1  \ begin{latexonly}
2  %<$separator%>
3  \ label{question:%<$label%>}
4  %<$prefix%>
5
6  \ begin{center}
7  \ begin{tabular}{llll}
8      $\bigcirc$ & TRUE & $\bigcirc$ & FALSE
9  \ end{tabular}
10
11  \noindent %<$sufix%>
12  \ end{center}
13  \ end{latexonly}
14
15  \ begin{htmlonly}
16  %<$separator%>
17  \ label{question:%<$label%>}
18  %<$prefix%>
19
20  \ begin{center}
21  \ begin{tabular}{llll}
22      \ href{$\bigcirc$}{answer:%<$label%>} & TRUE &
23      \ href{$\bigcirc$}{answer:%<$label%>} & FALSE
24  \ end{tabular}
25
26  \ noindent %<$sufix%>
27  \ end{center}
28  \ end{htmlonly}

```

El template se rellena en las líneas 87-94. En esa llamada se ejecuta el código Perl incrustado en el esqueleto y su resultado se inserta en la posición que ocupa en el texto.

Concatenación y Documentos HERE Cuando concatene sangre adecuadamente las concatenaciones:

```

my $usage = "Usage: $0 <file> [-full] [-o] [-beans]\n"
    . "Options:\n"
    . "    -full  : produce a full dump\n"
    . "    -o      : dump in octal\n"
    . "    -beans : source is Java\n"
    ;

```

ponga el punto al principio de la siguiente línea, no al final.

Pero cuando el número de líneas es grande es mejor usar un *here document* o *documento aquí*. Veamos un ejemplo:

```
print <<"EOI";
```

El programa se deberá ejecutar con:

```
$0 numfiles $opciones initialvalue
EOI
```

Para definir un “documento aquí” se escribe la etiqueta entrecomillada y precedida de << y sigue el texto que consituye el *here document* que se delimita por una línea en blanco que empieza por la etiqueta. Al documento aquí se le trata como una cadena de doble comilla si la etiqueta aparece en doble comilla y como de comilla simple si la etiqueta esta entre comillas simples. Observe que el punto y coma se escribe despues de la primera aparición de la etiqueta.

Un problema con el uso de los heredoc es que rompen la estructura normal del sangrado:

```
if ($usage_error) {
    warn <<'END_USAGE';
Usage: qdump <file> [-full] [-o] [-beans]
Options:
    -full   : produce a full dump
    -o      : dump in octal
    -beans  : source is Java
END_USAGE
}
```

Es mejor que cada heredoc se aise en una subrutina y se parametrize con las variables que van a ser interpoladas:

```
sub build_usage {
    my ($prog_name, $file_name) = @_;

    return <<"END_USAGE";
Usage: $prog_name $file_name [-full] [-o] [-beans]
Options:
    -full   : produce a full dump
    -o      : dump in octal
    -beans  : source is Java
END_USAGE
}
```

que mas tarde puede ser llamado con los valores de interpolación adecuados:

```
if ($usage_error) {
    warn build_usage($PROGRAM_NAME, $requested_file);
}
```

Véase el libro de Conway Perl Best Practices [?] para mas detalles sobre buenas prácticas de programación con heredocs.

Descarga de los Módulos Necesarios

- El módulo Grammar : <http://nereida.deioc.ull.es/~pl/perlexamples/Grammar-0.03.tar.gz>
- El módulo PL::FirstFollow : <http://nereida.deioc.ull.es/~pl/perlexamples/PL-FirstFollow-0.02.tar.gz>

4.7. Esquemas de Traducción

Definición 4.7.1. *Un esquema de traducción es una gramática independiente del contexto en la cual se han insertado fragmentos de código en las partes derechas de sus reglas de producción. Los fragmentos de código así insertados se denominan acciones semánticas. Dichos fragmentos actúan, calculan y modifican los atributos asociados con los nodos del árbol sintáctico. El orden en que se evalúan los fragmentos es el de un recorrido primero-profundo del árbol de análisis sintáctico.*

Obsérvese que, en general, para poder aplicar un esquema de traducción hay que construir el árbol sintáctico y después aplicar las acciones empujadas en las reglas en el orden de recorrido primero-profundo. Por supuesto, si la gramática es ambigua una frase podría tener dos árboles y la ejecución de las acciones para ellos podría dar lugar a diferentes resultados. Si se quiere evitar la multiplicidad de resultados (interpretaciones semánticas) es necesario precisar de que árbol sintáctico concreto se está hablando.

Por ejemplo, si en la regla $A \rightarrow \alpha\beta$ insertamos un fragmento de código:

$$A \rightarrow \alpha\{action\}\beta$$

La acción $\{action\}$ se ejecutará después de todas las acciones asociadas con el recorrido del subárbol de α y antes que todas las acciones asociadas con el recorrido del subárbol β .

El siguiente esquema de traducción recibe como entrada una expresión en infijo y produce como salida su traducción a postfijo para expresiones aritmeticas con sólo restas de números:

$$\begin{array}{ll} expr \rightarrow expr_1 - NUM & \{ \$expr\{TRA\} = \$expr[1]\{TRA\} "." ".$NUM\{VAL\}." - " \} \\ expr \rightarrow NUM & \{ \$expr\{TRA\} = $NUM\{VAL\} \} \end{array}$$

Las apariciones de variables sintácticas en una regla de producción se indexan como se ve en el ejemplo, para distinguir de que nodo del árbol de análisis estamos hablando. Cuando hablemos del atributo de un nodo utilizaremos una indexación tipo *hash*. Aquí VAL es un atributo de los nodos de tipo NUM denotando su valor numérico y para accederlo escribiremos \$NUM{VAL}. Análogamente \$expr{TRA} denota el atributo “traducción” de los nodos de tipo expr.

Ejercicio 4.7.1. *Muestre la secuencia de acciones a la que da lugar el esquema de traducción anterior para la frase 7 -5 -4.*

En este ejemplo, el cómputo del atributo \$expr{TRA} depende de los atributos en los nodos hijos, o lo que es lo mismo, depende de los atributos de los símbolos en la parte derecha de la regla de producción. Esto ocurre a menudo y motiva la siguiente definición:

Definición 4.7.2. *Un atributo tal que su valor en un nodo puede ser computado en términos de los atributos de los hijos del nodo se dice que es un atributo sintetizado.*

Ejemplo 4.7.1. *Un ejemplo de atributo heredado es el tipo de las variables en las declaraciones:*

$$\begin{array}{l} decl \rightarrow type \{ \$list\{T\} = \$type\{T\} \} list \\ type \rightarrow INT \{ \$type\{T\} = \$int \} \\ type \rightarrow STRING \{ \$type\{T\} = \$string \} \\ list \rightarrow ID , \{ \$ID\{T\} = \$list\{T\}; \$list_1\{T\} = \$list\{T\} \} list_1 \\ list \rightarrow ID \{ \$ID\{T\} = \$list\{T\} \} \end{array}$$

Definición 4.7.3. *Un atributo heredado es aquel cuyo valor se computa a partir de los valores de sus hermanos y de su padre.*

Ejercicio 4.7.2. *Escriba un esquema de traducción que convierta expresiones en infijo con los operadores +-*/() y números en expresiones en postfijo. Explique el significado de los atributos elegidos.*

4.8. Recursión por la Izquierda

Definición 4.8.1. Una gramática es recursiva por la izquierda cuando existe una derivación $A \xRightarrow{*} A\alpha$.

En particular, es recursiva por la izquierda si contiene una regla de producción de la forma $A \rightarrow A\alpha$. En este caso se dice que la recursión por la izquierda es directa.

Cuando la gramática es recursiva por la izquierda, el método de análisis recursivo descendente predictivo no funciona. En ese caso, el procedimiento A asociado con A ciclaría para siempre sin llegar a consumir ningún terminal.

4.8.1. Eliminación de la Recursión por la Izquierda en la Gramática

Es posible modificar la gramática para eliminar la recursión por la izquierda. En este apartado nos limitaremos al caso de recursión por la izquierda directa. La generalización al caso de recursión por la izquierda no-directa se reduce a la iteración de la solución propuesta para el caso directo.

Consideremos una variable A con dos producciones:

$$A \rightarrow A\alpha \quad | \quad \beta$$

donde $\alpha, \beta \in (V \cup \Sigma)^*$ no comienzan por A. Estas dos producciones pueden ser sustituidas por:

$$\begin{aligned} A &\rightarrow \beta R \\ R &\rightarrow \alpha R \quad | \quad \epsilon \end{aligned}$$

eliminando así la recursión por la izquierda.

Definición 4.8.2. La producción $R \rightarrow \alpha R$ se dice recursiva por la derecha.

Las producciones recursivas por la derecha dan lugar a árboles que se hunden hacia la derecha. Es más difícil traducir desde esta clase de árboles operadores como el menos, que son asociativos a izquierdas.

Ejercicio 4.8.1. Elimine la recursión por la izquierda de la gramática

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} - \text{NUM} \\ \text{expr} &\rightarrow \text{NUM} \end{aligned}$$

Ejercicio 4.8.2. ¿Que hay de erróneo en este esquema de traducción?

$$\begin{aligned} \text{expr} &\rightarrow \text{NUM} - \text{expr}_1 & \{ \text{\$expr\{T\}} = \text{\$NUM\{VAL\}}." ".\text{\$expr[1]\{T\}}." - " \} \\ \text{expr} &\rightarrow \text{NUM} & \{ \text{\$expr\{T\}} = \text{\$NUM\{VAL\}} \} \end{aligned}$$

Ejercicio 4.8.3. Dado el esquema de traducción:

$$\begin{aligned} e &\rightarrow \text{NUM } r & \{ \text{\$e\{TRA\}} = \text{\$NUM\{VAL\}}." ".\text{\$r\{TRA\}} \} \\ r &\rightarrow -e & \{ \text{\$r\{TRA\}} = \text{\$e\{TRA\}}." - " \} \\ r &\rightarrow \epsilon & \{ \text{\$r\{TRA\}} = "" \} \end{aligned}$$

¿Cuál es el lenguaje generado por la gramática? ¿Puede el lenguaje ser analizado por un APDR? ¿Cual es la traducción de 4-5-6? ¿Es un esquema de traducción adecuado para traducir de infijo a postfijo? ¿Cuál es la traducción si cambiamos el anterior esquema por este otro?:

$$\begin{aligned} e &\rightarrow \text{NUM } r & \{ \text{\$e\{TRA\}} = \text{\$NUM\{VAL\}}." ".\text{\$r\{TRA\}} \} \\ r &\rightarrow -e & \{ \text{\$r\{TRA\}} = " - ".\text{\$e\{TRA\}} \} \\ r &\rightarrow \epsilon & \{ \text{\$r\{TRA\}} = "" \} \end{aligned}$$

4.8.2. Eliminación de la Recursión por la Izquierda en un Esquema de Traducción

La eliminación de la recursión por la izquierda es sólo un paso: debe ser extendida a esquemas de traducción, de manera que no sólo se preserve el lenguaje sino la secuencia de acciones. Supongamos que tenemos un esquema de traducción de la forma:

$$\begin{aligned} A &\rightarrow A\alpha \quad \{ \text{alpha_action} \} \\ A &\rightarrow A\beta \quad \{ \text{beta_action} \} \\ A &\rightarrow \gamma \quad \{ \text{gamma_action} \} \end{aligned}$$

para una sentencia como $\gamma\beta\alpha$ la secuencia de acciones será:

$$\text{gamma_action} \quad \text{beta_action} \quad \text{alpha_action}$$

¿Cómo construir un esquema de traducción para la gramática resultante de eliminar la recursión por la izquierda que ejecute las acciones asociadas en el mismo orden?. Supongamos para simplificar, que las acciones no dependen de atributos ni computan atributos, sino que actúan sobre variables globales. En tal caso, la siguiente ubicación de las acciones da lugar a que se ejecuten en el mismo orden:

$$\begin{aligned} A &\rightarrow \gamma \{ \text{gamma_action} \} R \\ R &\rightarrow \beta \{ \text{beta_action} \} R \\ R &\rightarrow \alpha \{ \text{alpha_action} \} R \\ R &\rightarrow \epsilon \end{aligned}$$

Si hay atributos en juego, la estrategia para construir un esquema de traducción equivalente para la gramática resultante de eliminar la recursividad por la izquierda se complica. Consideremos de nuevo el esquema de traducción de infijo a postfijo de expresiones aritméticas de restas:

$$\begin{aligned} \text{expr} &\rightarrow \text{expr}_1 - \text{NUM} \quad \{ \$\text{expr}\{\text{T}\} = \$\text{expr}[1]\{\text{T}\} \cdot " \cdot \$\text{NUM}\{\text{VAL}\} \cdot " - " \} \\ \text{expr} &\rightarrow \text{NUM} \quad \{ \$\text{expr}\{\text{T}\} = \$\text{NUM}\{\text{VAL}\} \} \end{aligned}$$

En este caso introducimos un atributo H para los nodos de la clase r el cuál acumula la traducción a postfijo hasta el momento. Observe como este atributo se computa en un nodo r a partir del correspondiente atributo del el padre y/o de los hermanos del nodo:

$$\begin{aligned} \text{expr} &\rightarrow \text{NUM} \{ \$r\{\text{H}\} = \$\text{NUM}\{\text{VAL}\} \} r \{ \$\text{expr}\{\text{T}\} = \$r\{\text{T}\} \} \\ r &\rightarrow -\text{NUM} \{ \$r_1\{\text{H}\} = \$r\{\text{H}\} \cdot " \cdot \$\text{NUM}\{\text{VAL}\} \cdot " - " \} r_1 \{ \$r\{\text{T}\} = \$r_1\{\text{T}\} \} \\ r &\rightarrow \epsilon \{ \$r\{\text{T}\} = \$r\{\text{H}\} \} \end{aligned}$$

El atributo H es un ejemplo de atributo heredado.

4.8.3. Ejercicio

Calcule los valores de los atributos cuando se aplica el esquema de traducción anterior a la frase $4 - 5 - 7$.

4.8.4. Convirtiendo el Esquema en un Analizador Predictivo

A partir del esquema propuesto, que se basa en una fase de descenso con un atributo heredado y una de ascenso con un atributo sintetizado:

$$\begin{aligned} \text{expr} &\rightarrow \text{NUM} \{ \$r\{\text{H}\} = \$\text{NUM}\{\text{VAL}\} \} r \{ \$\text{expr}\{\text{T}\} = \$r\{\text{T}\} \} \\ r &\rightarrow -\text{NUM} \{ \$r_1\{\text{H}\} = \$r\{\text{H}\} \cdot " \cdot \$\text{NUM}\{\text{VAL}\} \cdot " - " \} r_1 \{ \$r\{\text{T}\} = \$r_1\{\text{T}\} \} \\ r &\rightarrow \epsilon \{ \$r\{\text{T}\} = \$r\{\text{H}\} \} \end{aligned}$$

es posible construir un APDR que ejecuta las acciones semánticas en los puntos indicados por el esquema de traducción. El atributo heredado se convierte en un parámetro de entrada a la subrutina asociada con la variable sintáctica:

```
sub expression() {
  my $r = $value." "; #accion intermedia
  match('NUM');
  return rest($r); # accion final $expr{T} = $r{T}
}

sub rest($) {
  my $v = shift;

  if ($lookahead eq '-') {
    match('-');
    my $r = "$v $value -"; # accion intermedia
    match('NUM');
    return rest($r); # accion final $r{t} = $r_1{t}
  }
  elsif ($lookahead ne 'EOI') {
    error("Se esperaba un operador");
  }
  else { return $v; } # r -> epsilon { $r{t} = $r{h} }
}
```

4.8.5. Ejercicio

Generalize la estrategia anterior para eliminar la recursividad por la izquierda al siguiente esquema de traducción genérico recursivo por la izquierda y con un atributo sintetizado A^s :

$$\begin{array}{ll} A \rightarrow A_1 X_1 X_2 X_3 & \{A^s = f_X(A_1^s, X_1^s, X_2^s, X_3^s)\} \\ A \rightarrow A_1 Y_1 Y_2 Y_3 & \{A^s = f_Y(A_1^s, Y_1^s, Y_2^s, Y_3^s)\} \\ A \rightarrow Z_1 Z_2 Z_3 & \{A^s = f_Z(Z_1^s, Z_2^s, Z_3^s)\} \end{array}$$

donde f_X, f_Y y f_Z son funciones cualesquiera.

4.8.6. Práctica: Eliminación de la Recursividad por la Izquierda

En esta práctica vamos a extender las fases de análisis léxico y sintáctico del compilador del lenguaje Tutu cuya gramática se definió en el ejercicio 4.5.1 con expresiones que incluyen diferencias y divisiones. Además construiremos una representación del árbol sintáctico concreto. Para ello consideremos el siguiente esquema de traducción recursivo por la izquierda (en concreto las reglas recursivas por la izquierda son las 10, 11, 13 y 14):

0	$p \rightarrow ds\ ss$	$\{ \$p\{t\} = \{ n \Rightarrow 0, ds \Rightarrow \$ds\{t\}, ss \Rightarrow \$ss\{t\} \} \}$
1	$p \rightarrow ss$	$\{ \$p\{t\} = \{ n \Rightarrow 1, ss \Rightarrow \$ss\{t\} \} \}$
2	$ds \rightarrow d\ ';' \ ds$	$\{ \$ds\{t\} = \{ n \Rightarrow 2, d \Rightarrow \$d\{t\}, ; \Rightarrow ';' , ds \Rightarrow \$ds\{t\} \} \}$
3	$ds \rightarrow d\ ';' \$	$\{ \$ds\{t\} = \{ n \Rightarrow 3, d \Rightarrow \$d\{t\} ; \Rightarrow ';' \} \}$
4	$d \rightarrow INT\ il$	$\{ \$d\{t\} = \{ n \Rightarrow 4, INT \Rightarrow 'INT', il \Rightarrow \$il\{t\} \} \}$
5	$d \rightarrow STRING\ il$	$\{ \$d\{t\} = \{ n \Rightarrow 5, STRING \Rightarrow 'STRING', il \Rightarrow \$il\{t\} \} \}$
6	$ss \rightarrow s\ ';' \ ss$	$\{ \$ss\{t\} = \{ n \Rightarrow 6, s \Rightarrow \$s\{t\}, ; \Rightarrow ';' , ss \Rightarrow \$ss\{t\} \} \}$
7	$ss \rightarrow s$	$\{ \$ss\{t\} = \{ n \Rightarrow 7, s \Rightarrow \$s\{t\} \} \}$
8	$s \rightarrow ID\ '=' \ e$	$\{ \$s\{t\} = \{ n \Rightarrow 8, ID \Rightarrow \$ID\{v\}, = \Rightarrow '=' , e \Rightarrow \$e\{t\} \} \}$
9	$s \rightarrow P\ e$	$\{ \$s\{t\} = \{ n \Rightarrow 9, P \Rightarrow 'P', e \Rightarrow \$e\{t\} \} \}$
10	$e \rightarrow e1\ '+' \ t$	$\{ \$e\{t\} = \{ n \Rightarrow 10, e \Rightarrow \$e1\{t\}, + \Rightarrow '+' , t \Rightarrow \$t\{t\} \} \}$
11	$e \rightarrow e1\ '-' \ t$	$\{ \$e\{t\} = \{ n \Rightarrow 11, e \Rightarrow \$e1\{t\}, - \Rightarrow '-' , t \Rightarrow \$t\{t\} \} \}$
12	$e \rightarrow t$	$\{ \$e\{t\} = \{ n \Rightarrow 12, t \Rightarrow \$t\{t\} \} \}$
13	$t \rightarrow t1\ '*' \ f$	$\{ \$t\{t\} = \{ n \Rightarrow 13, t \Rightarrow \$t1\{t\}, * \Rightarrow '*' , f \Rightarrow \$f\{t\} \} \}$
14	$t \rightarrow t\ '/' \ f$	$\{ \$t\{t\} = \{ n \Rightarrow 14, t \Rightarrow \$t1\{t\}, / \Rightarrow '/' , f \Rightarrow \$f\{t\} \} \}$
15	$t \rightarrow f$	$\{ \$t\{t\} = \{ n \Rightarrow 15, f \Rightarrow \$f\{t\} \} \}$
16	$f \rightarrow '(' \ e \ ')'$	$\{ \$f\{t\} = \{ n \Rightarrow 16, (\Rightarrow '(' , e \Rightarrow \$e\{t\},) \Rightarrow ')' \} \}$
17	$f \rightarrow ID$	$\{ \$f\{t\} = \{ n \Rightarrow 17, ID \Rightarrow \$ID\{v\} \} \}$
18	$f \rightarrow NUM$	$\{ \$f\{t\} = \{ n \Rightarrow 18, NUM \Rightarrow \$NUM\{v\} \} \}$
19	$f \rightarrow STR$	$\{ \$f\{t\} = \{ n \Rightarrow 19, STR \Rightarrow \$STR\{v\} \} \}$
20	$il \rightarrow ID\ ',' \ il$	$\{ \$il\{t\} = \{ n \Rightarrow 20, ID \Rightarrow \$ID\{v\}, ',' \Rightarrow ',' , il \Rightarrow \$il\{t\} \} \}$
21	$il \rightarrow ID$	$\{ \$il\{t\} = \{ n \Rightarrow 21, ID \Rightarrow \$ID\{v\} \} \}$
22	$s \rightarrow \epsilon$	$\{ \$s\{t\} = \{ n \Rightarrow 22, s \Rightarrow '' \} \}$

Por razones de espacio hemos abreviado los nombres de las variables. El atributo t (por *tree*) es una referencia a un hash. La entrada n contiene el número de la regla en juego. Hay una entrada por símbolo en la parte derecha. El atributo v de ID es la cadena asociada con el identificador. El atributo v de NUM es el valor numérico asociado con el terminal. Se trata de, siguiendo la metodología explicada en la sección anterior, construir un analizador descendente predictivo recursivo que sea equivalente al esquema anterior. Elimine la recursión por la izquierda. Traslade las acciones a los lugares convenientes en el nuevo esquema e introduzca los atributos heredados que sean necesarios. Genere pruebas siguiendo la metodología explicada en la sección 4.4.1. ¡Note que el árbol que debe producir es el de la gramática inicial, ¡No el de la gramática transformada!

4.9. Árbol de Análisis Abstracto

4.9.1. Lenguajes Árbol y Gramáticas Árbol

Un *árbol de análisis abstracto* (denotado AAA, en inglés *abstract syntax tree* o *AST*) porta la misma información que el árbol de análisis sintáctico pero de forma mas condensada, eliminándose terminales y producciones que no aportan información.

Definición 4.9.1. Un alfabeto con función de aridad es un par (Σ, ρ) donde Σ es un conjunto finito y ρ es una función $\rho : \Sigma \rightarrow \mathbb{N}_0$, denominada función de aridad. Denotamos por $\Sigma_k = \{a \in \Sigma : \rho(a) = k\}$.

Definimos el lenguaje árbol homogéneo $B(\Sigma)$ sobre Σ inductivamente:

- Todos los elementos de aridad 0 están en $B(\Sigma)$: $a \in \Sigma_0$ implica $a \in B(\Sigma)$
- Si $b_1, \dots, b_k \in B(\Sigma)$ y $f \in \Sigma_k$ es un elemento k -ario, entonces $f(b_1, \dots, b_k) \in B(\Sigma)$

Los elementos de $B(\Sigma)$ se llaman árboles o términos.

Ejemplo 4.9.1. Sea $\Sigma = \{A, CONS, NIL\}$ con $\rho(A) = \rho(NIL) = 0, \rho(CONS) = 2$. Entonces $B(\Sigma) = \{A, NIL, CONS(A, NIL), CONS(NIL, A), CONS(A, A), CONS(NIL, NIL), \dots\}$

Ejemplo 4.9.2. Una versión simplificada del alfabeto con aridad en el que estan basados los árboles construidos por el compilador de Tutu es:

$\Sigma = \{ID, NUM, LEFTVALUE, STR, PLUS, TIMES, ASSIGN, PRINT\}$

$\rho(ID) = \rho(NUM) = \rho(LEFTVALUE) = \rho(STR) = 0$

$\rho(PRINT) = 1$

$\rho(PLUS) = \rho(TIMES) = \rho(ASSIGN) = 2$.

Observe que los elementos en $B(\Sigma)$ no necesariamente son árboles “correctos”. Por ejemplo, el árbol $ASSIGN(NUM, PRINT(ID))$ es un elemento de $B(\Sigma)$.

Definición 4.9.2. Una gramática árbol regular es una cuadrupla $((\Sigma, \rho), N, P, S)$, donde:

- (Σ, ρ) es un alfabeto con aricidad $\rho : \Sigma \rightarrow \mathbb{N}$
- N es un conjunto finito de variables sintácticas o no terminales
- P es un conjunto finito de reglas de producción de la forma $X \rightarrow s$ con $X \in N$ y $s \in B(\Sigma \cup N)$
- $S \in N$ es la variable o símbolo de arranque

Definición 4.9.3. Dada una gramática (Σ, N, P, S) , se dice que un árbol $t \in B(\Sigma \cup N)$ es del tipo (X_1, \dots, X_k) si el j -ésimo noterminal, contando desde la izquierda, que aparece en t es $X_j \in N$.

Si $p = X \rightarrow s$ es una producción y s es de tipo (X_1, \dots, X_n) , diremos que la producción p es de tipo $(X_1, \dots, X_n) \rightarrow X$.

Definición 4.9.4. Consideremos un árbol $t \in B(\Sigma \cup N)$ que sea del tipo (X_1, \dots, X_n) , esto es las variables sintácticas en el árbol leídas de izquierda a derecha son (X_1, \dots, X_n) .

- Si $X_i \rightarrow s_i \in P$ para algún i , entonces decimos que el árbol t deriva en un paso en el árbol t' resultante de sustituir el nodo X_i por el árbol s_i y escribiremos $t \Rightarrow t'$. Esto es, $t' = t\{X_i/s_i\}$
- Todo árbol deriva en cero pasos en si mismo $t \xRightarrow{0} t$.
- Decimos que un árbol t deriva en n pasos en el árbol t' y escribimos $t \xRightarrow{n} t'$ si t deriva en un paso en un árbol t'' el cuál deriva en $n - 1$ pasos en t' . En general, si t deriva en un cierto número de pasos en t' escribiremos $t \xRightarrow{*} t'$.

Definición 4.9.5. Se define el lenguaje árbol generado por una gramática $G = (\Sigma, N, P, S)$ como el lenguaje $L(G) = \{t \in B(\Sigma) : \exists S \xRightarrow{*} t\}$.

Ejemplo 4.9.3. Sea $G = (\Sigma, V, P, S)$ con $\Sigma = \{A, CONS, NIL\}$ y $\rho(A) = \rho(NIL) = 0, \rho(CONS) = 2$ y sea $V = \{E, L\}$. El conjunto de producciones P es:

$$P_1 = \{L \rightarrow NIL, L \rightarrow CONS(E, L), E \rightarrow a\}$$

La producción $L \rightarrow CONS(E, L)$ es del tipo $(E, L) \rightarrow L$.

Informalmente, el lenguaje generado por G se obtiene realizando sustituciones sucesivas (derivando) desde el símbolo de arranque hasta producir un árbol cuyos nodos estén etiquetados con elementos de Σ . Debería ser claro que, en este ejemplo, $L(G)$ es el conjunto de las listas en A , incluyendo la lista vacía:

$$L(G) = \{NIL, CONS(A, NIL), CONS(A, CONS(A, NIL)), \dots\}$$

Ejercicio 4.9.1. Construya una derivación para el árbol $CONS(A, CONS(A, NIL))$. ¿De que tipo es el árbol $CONS(E, CONS(A, CONS(E, L)))$?

Cuando hablamos del AAA producido por un analizador sintáctico, estamos en realidad hablando de un lenguaje árbol cuya definición precisa debe hacerse a través de una gramática árbol regular. Mediante las gramáticas árbol regulares disponemos de un mecanismo para describir formalmente el lenguaje de los AAA que producirá el analizador sintáctico para las sentencias Tutu.

Ejemplo 4.9.4. Sea $G = (\Sigma, V, P, S)$ con

$$\begin{aligned}\Sigma &= \{ID, NUM, LEFTVALUE, STR, PLUS, TIMES, ASSIGN, PRINT\} \\ \rho(ID) &= \rho(NUM) = \rho(LEFTVALUE) = \rho(STR) = 0 \\ \rho(PRINT) &= 1 \\ \rho(PLUS) &= \rho(TIMES) = \rho(ASSIGN) = 2 \\ V &= \{st, expr\}\end{aligned}$$

y las producciones:

$$P = \left\{ \begin{array}{ll} st & \rightarrow ASSIGN(LEFTVALUE, expr) \\ st & \rightarrow PRINT(expr) \\ expr & \rightarrow PLUS(expr, expr) \\ expr & \rightarrow TIMES(expr, expr) \\ expr & \rightarrow NUM \\ expr & \rightarrow ID \\ expr & \rightarrow STR \end{array} \right\}$$

Entonces el lenguaje $L(G)$ contiene árboles como el siguiente:

$$\begin{array}{l} ASSIGN \ (\\ \quad LEFTVALUE, \\ \quad PLUS \ (\\ \qquad ID, \\ \qquad TIMES \ (\\ \qquad \quad NUM, \\ \qquad \quad ID \\ \qquad \) \\ \) \\ \) \end{array}$$

El cual podría corresponderse con una sentencia como $a = b + 4 * c$.

El lenguaje de árboles descrito por esta gramática árbol es el lenguaje de los AAA de las sentencias de Tutu.

Ejercicio 4.9.2. Redefina el concepto de árbol de análisis concreto dado en la definición 8.1.7 utilizando el concepto de gramática árbol. Con mas precisión, dada una gramática $G = (\Sigma, V, P, S)$ defina una gramática árbol $T = (\Omega, N, R, U)$ tal que $L(T)$ sea el lenguaje de los árboles concretos de G . Puesto que las partes derechas de las reglas de producción de P pueden ser de distinta longitud, existe un problema con la aricidad de los elementos de Ω . Discuta posibles soluciones.

Ejercicio 4.9.3. ¿Cómo son los árboles sintácticos en las derivaciones árbol? Dibuje varios árboles sintácticos para las gramáticas introducidas en los ejemplos 4.9.3 y 4.9.4.

Intente dar una definición formal del concepto de árbol de análisis sintáctico asociado con una derivación en una gramática árbol

Definición 4.9.6. La notación de Dewey es una forma de especificar los subárboles de un árbol $t \in B(\Sigma)$. La notación sigue el mismo esquema que la numeración de secciones en un texto: es una palabra formada por números separados por puntos. Así $t/2.1.3$ denota al tercer hijo del primer hijo del segundo hijo del árbol t . La definición formal sería:

- $t/\epsilon = t$
- Si $t = a(t_1, \dots, t_k)$ y $j \in \{1 \dots k\}$ y n es una cadena de números y puntos, se define inductivamente el subárbol $t/j.n$ como el subárbol n -ésimo del j -ésimo subárbol de t . Esto es: $t/j.n = t_j/n$

Ejercicio 4.9.4. Sea el árbol:

$$t = \text{ASSIGN} \left(\begin{array}{l} \text{LEFTVALUE}, \\ \text{PLUS} \end{array} \left(\begin{array}{l} \text{ID}, \\ \text{TIMES} \left(\begin{array}{l} \text{NUM}, \\ \text{ID} \end{array} \right) \end{array} \right) \right)$$

Calcule los subárboles t/ϵ , $t/2.2.1$, $t/2.1$ y $t/2.1.2$.

4.9.2. Realización del AAA para Tutu en Perl

En la sección 4.6.1 nos limitamos a realizar un recorrido del árbol de análisis sintáctico concreto. En esta sección construimos un *árbol de análisis sintáctico abstracto*. Este proceso puede verse como la traducción desde el lenguaje de árboles concretos hasta el lenguaje de árboles abstractos.

Definición 4.9.7. La gramática árbol extendida que especifica los árboles AAA para el compilador de Tutu es esta:

1	$prog$	$\rightarrow PROGRAM(decls, sts)$
2	$decls$	$\rightarrow list\ decl$
3	sts	$\rightarrow list\ st$
4	$decl$	$\rightarrow INT(idlist)$
5	$decl$	$\rightarrow STRING(idlist)$
6	$idlist$	$\rightarrow list\ SIMPLEID$
7	st	$\rightarrow ASSIGN(LEFTVALUE, expr)$
8	st	$\rightarrow PRINT(expr)$
9	$expr$	$\rightarrow PLUS(expr, expr)$
10	$expr$	$\rightarrow TIMES(expr, expr)$
11	$expr$	$\rightarrow NUM$
12	$expr$	$\rightarrow ID$
13	$expr$	$\rightarrow STR$

Hemos extendido el concepto de gramática árbol con el concepto de *lista de no terminales*. A la hora de construir las estructuras de datos las listas de variables se van a traducir por listas de árboles.

Por ejemplo, un árbol abstracto para el programa

```
int a,b;
string c, d;
a = 4;
p a
```

Sería de la forma:

```
PROGRAM(
  DECLS[INT[ID, ID], STRING[ID, ID]],
  STS[ASSIGN(LEFTVALUE, NUM), PRINT(ID)]
)
```

Donde los corchetes indican listas y los paréntesis tuplas.

Para llevar a cabo la traducción deberemos tomar decisiones sobre que forma de representación nos conviene. Cada nodo del AAA va a ser un objeto y la clase indicará si es un nodo suma, producto, una declaración, una asignación, etc.

Cada nodo del árbol AAA va a ser un objeto. De este modo el acceso a los atributos del nodo se hará a través de los métodos asociados. Además, el procedimiento de traducción al lenguaje objetivo depende del tipo de nodo. Así por ejemplo, el método *traducción* es diferente para un nodo de tipo *PLUS* que para otro de tipo *ASSIGN*.

Resumamos antes de entrar en detalle, la forma de manejar los objetos en Perl:

- Para crear una *clase* se construye un “package”:

```
package NUM;
```

- Para crea un *método* se escribe una subrutina:

```
package NUM;
sub incr { my $self = shift; $self->{VAL}++ }
```

el primer argumento de un método suele ser la referencia al objeto en cuestión.

- Para crear un *objeto*, se bendice (“bless”) una referencia. Los objetos Perl son datos normales como hashes y arrays que han sido “bendecidos” en un paquete. Por ejemplo:

```
my $a = bless {VAL => 4}, 'NUM';
```

crea un objeto referenciado por **\$a** que pertenece a la clase NUM. Los métodos del objeto son las subrutinas que aparecen en el **package NUM**.

- Para referirse a un método de un objeto se usa la sintáxis “flecha”:

```
$a->incr;
```

Cuando se usa la sintáxis flecha, el primer argumento de la rutina es la referencia al objeto, esto es, la llamada anterior es equivalente a **NUM::incr(\$a)**

- Constructores: En Perl son rutinas que retornan una referencia a un objeto recién creado e inicializado

```
sub new { my ($class, $value) = @_; return bless {VAL => $value}, $class; }
```

Normalmente se llaman usando la sintáxis flecha, pero a la izquierda de la flecha va el nombre de la clase. Por ejemplo:

```
my $a = NUM->new(4)
```

En este caso, el primer argumento es el nombre de la clase. La llamada anterior es equivalente a **NUM::new('NUM', 4)**

Volviendo a nuestro problema de crear el AAA, para crear los objetos de las diferentes clases de nodos usaremos el módulo **Class::MakeMethods::Emulator::MethodMaker** (véase la línea 9):

```
1 package PL::Syntax::Analysis;
2
3 use 5.006;
4 use strict;
5 use warnings;
6 use Data::Dumper;
7 use IO::File;
8 use Class::MakeMethods::Emulator::MethodMaker '-sugar';
9
```

```

10 require Exporter;
11 our @ISA = qw(Exporter);
12 our @EXPORT = qw( );
13 our $VERSION = '0.02';
14
15 #####
16
17 # Grammar:
18 # P : DD L      | L
19 # DD: D ';' DD  | D ';'
20 # D : int I     | string I
21 # L : S         | S ; L
22 # S : ID '=' E  | p E | epsilon
23 # E : T '+' E   | T
24 # T : F '*' T   | F
25 # F : '(' E ')' | id | num | str
26 # I : id ',' I  | id

```

Hemos aislado la fase de análisis sintáctica en un módulo aparte denominado `PL::Syntax::Analysis`. La dependencia se actualiza en `Makefile.PL`:

```

PL0506/04sintactico/PL-Tutu$ cat -n Makefile.PL
 1 use 5.008004;
 2 use ExtUtils::MakeMaker;
 3 WriteMakefile(
 4     NAME          => 'PL::Tutu',
 5     VERSION_FROM  => 'lib/PL/Tutu.pm', # finds $VERSION
 6     PREREQ_PM     => {Class::MakeMethods::Emulator::MethodMaker => 0},1
 7     EXE_FILES     => [ 'scripts/tutu.pl', 'scripts/tutu' ],
 8     ($] >= 5.005 ? ## Add these new keywords supported since 5.005
 9         (ABSTRACT_FROM => 'lib/PL/Tutu.pm', # retrieve abstract from module
10         AUTHOR        => 'Casiano Rodriguez Leon <casiano@ull.es>') : ()),
11 );

```

Se actualiza también `MANIFEST`:

```

$ cat -n MANIFEST
 1 Changes
 2 lib/PL/Error.pm
 3 lib/PL/Lexical/Analysis.pm
 4 lib/PL/Syntax/Analysis.pm
 5 lib/PL/Tutu.pm
 6 Makefile.PL
 7 MANIFEST
 8 MANIFEST.SKIP
 9 README
10 scripts/test01.tutu
11 scripts/tutu
12 scripts/tutu.pl
13 t/01Lexical.t

```

Ahora compile llama a `Syntax::Analysis::parser` pasándole como argumento la lista de terminales `@tokens`. La función `parser` devuelve el AAA:

```

04sintactico/PL-Tutu/lib/PL$ sed -ne '/sub compile\>/,/^\}/p' Tutu.pm | cat -n
 1 sub compile {

```

```

2   my ($input) = @_;
3   #my %symbol_table = ();
4   #my $data = ""; # Contiene todas las cadenas en el programa fuente
5   my $target = ""; # target code
6   my @tokens = ();
7   my $tree = undef; # abstract syntax tree
8   #my $global_address = 0;
9
10
11  #####lexical analysis
12  @tokens = scanner($input);
13  #print "@tokens\n";
14
15  #####syntax (and semantic) analysis
16  $tree = &Syntax::Analysis::parser(@tokens);
17  print Dumper($tree);
18
19  #####machine independent optimizations
20  &Machine::Independent::Optimization::Optimize;
21
22  #####code generation
23  &Code::Generation::code_generator;
24
25  #####peephole optimization
26  &Peephole::Optimization::transform($target);
27
28  return \$target;
29 }

```

El módulo `Class::MakeMethods::Emulator::MethodMaker` permite crear constructores y métodos de acceso. El módulo no viene con la distribución de Perl, así que, en general, deberá descargarlo desde CPAN e instalarlo. Así definimos que existe una clase de nodos `TYPE` que nuestro AAA va a tener:

```

package TYPE;
make methods
    get_set      => [ 'NAME', 'LENGTH' ],
    new_hash_init => 'new';

```

El uso de los argumentos `get_set => ['NAME', 'LENGTH']` hace que se cree un objeto de tipo hash con claves `'NAME'` y `'LENGTH'` así como métodos `NAME` y `LENGTH` que cuando se llaman con un argumento devuelven el valor y cuando se llaman con dos argumentos modifican el valor correspondiente. La clave `get_set` produce métodos de acceso y modificación de los atributos del objeto que tienen la forma:

```

sub TYPE::NAME {
    my ($self, $new) = @_;
    defined($new) and $self->{NAME} = $new;
    return $self->{NAME};
}

```

Así mismo el uso de `new_hash_init => 'new'` genera un constructor cuyo nombre será `'new'` y que cuando es llamado inicializará el objeto con los argumentos con nombre especificados en la llamada. El constructor construido (vaya retruécano) cuando se usa la clave `new_hash_init` tiene el siguiente aspecto:

```

sub TYPE::new {
    my ($class, %args) = @_;
    my $self = {};

    bless $self, $class;
    foreach my $attribute (keys %args) {
        $self->$attribute($args{$attribute});
    }
    return $self;
}

```

Ahora podemos crear objetos de la clase TYPE haciendo:

```

my $int_type = TYPE->new(NAME => 'INTEGER', LENGTH => 1);
my $string_type = TYPE->new(NAME => 'STRING', LENGTH => 2);
my $err_type = TYPE->new(NAME => 'ERROR', LENGTH => 0);

```

Cada uno de estos objetos es un hash con las correspondientes claves para el nombre y el tipo.

Otros tipos de nodos del AAA son:

```

package PROGRAM; # raíz del AAA
make methods
    get_set      => [ 'DECLS', 'STS' ],
    new_hash_init => 'new';

package STRING; # tipo
make methods
    get_set      => [ 'TYPE', 'IDLIST' ],
    new_hash_init => 'new';

package INT; # tipo
make methods
    get_set      => [ 'TYPE', 'IDLIST' ],
    new_hash_init => 'new';

package ASSIGN; #sentencia
make methods
    get_set      => [ 'LEFT', 'RIGHT' ],
    new_hash_init => 'new';

package PRINT; #sentencia
make methods
    get_set      => [ 'EXPRESSION' ],
    new_hash_init => 'new';

package NUM; # para los números
make methods
    get_set      => [ 'VAL', 'TYPE' ],
    new_hash_init => 'new';

package ID; # Nodos identificador. Parte derecha
make methods
    get_set      => [ 'VAL', 'TYPE' ],
    new_hash_init => 'new';

```

```

package STR; # Clase para las constantes cadena
make methods
  get_set      => [ 'OFFSET', 'LENGTH', 'TYPE' ],
  new_hash_init => 'new';

package PLUS; # Nodo suma
make methods
  get_set      => [ 'LEFT', 'RIGHT', 'TYPE' ],
  new_hash_init => 'new';

package TIMES;
make methods
  get_set      => [ 'LEFT', 'RIGHT', 'TYPE' ],
  new_hash_init => 'new';

package LEFTVALUE; # Identificador en la parte izquierda
make methods      # de una asignación
  get_set      => [ 'VAL', 'TYPE' ],
  new_hash_init => 'new';

```

Hemos extendido el concepto de gramática árbol con el concepto de *lista de no terminales*. A la hora de construir las estructuras de datos las listas de variables se van a traducir por listas de árboles. Los tipos de nodos (*ASSIGN*, *PRINT*, ...) se traducen en nombres de clases. Hemos hecho una excepción con *SIMPLEID* el cual es simplemente una variable cadena conteniendo el identificador correspondiente.

El siguiente esquema de traducción resume la idea para una gramática simplificada: cada vez que encontremos un nodo en el árbol sintáctico concreto con una operación crearemos un nodo en el AAA cuya clase viene definida por el tipo de operación. Para los terminales creamos igualmente nodos indicando de que clase de terminal se trata. El atributo nodo lo denotaremos por *n*:

$$\begin{array}{ll}
 e \rightarrow e_1 + f & \{ \$e\{n\} = \text{PLUS} \rightarrow \text{new}(\text{LEFT} \Rightarrow \$e_1\{n\}, \text{RIGHT} \Rightarrow \$f\{n\}) \} \\
 f \rightarrow \text{NUM} & \{ \$f\{n\} = \text{NUM} \rightarrow \text{new}(\text{VAL} \Rightarrow \$\text{NUM}\{\text{VAL}\}) \} \\
 f \rightarrow \text{ID} & \{ \$f\{n\} = \text{ID} \rightarrow \text{new}(\text{VAL} \Rightarrow \$\text{ID}\{\text{VAL}\}) \}
 \end{array}$$

La estructura de cada rutina sigue siendo la misma, sólo que ampliada con las acciones para la construcción de los correspondientes nodos. Veamos por ejemplo, como modificamos la subrutina factor:

```

sub factor() {
  my ($e, $id, $str, $num);

  if ($lookahead eq 'NUM') {
    $num = $value;
    match('NUM');
    return NUM->new(VAL => $num, TYPE => $int_type);
  }
  elsif ($lookahead eq 'ID') {
    $id = $value;
    match('ID');
    return ID->new( VAL => $id, TYPE => undef);
  }
  elsif ($lookahead eq 'STR') {
    $str = $value;
    match('STR');
    return STR->new(OFFSET => undef, LENGTH => undef, TYPE => $string_type);
  }
}

```



```

}
elseif ($lookahead eq '(') {
    match('(');
    $e = expression;
    match(')');
    return $e;
}
else {
    Error::fatal("Se esperaba (, NUM o ID");
}
}

```

4.9.3. AAA: Otros tipos de nodos

Hemos optado por que las rutinas asociadas a variables sintácticas que describen listas de subcategorías devuelvan las correspondientes listas de nodos. Teníamos tres variables tipo lista. Las reglas para las listas eran:

Gramática de los Árboles de Tutu	Gramática del lenguaje Tutu
$decls \rightarrow list\ decl$	$declarations \rightarrow declaration\ ';' declarations \mid declaration\ ';'$
$sts \rightarrow list\ st$	$statements \rightarrow statement\ ';' statements \mid statement$
$idlist \rightarrow list\ SIMPLEID$	$idlist \rightarrow ID\ ',' idlist \mid ID$

En este caso las subrutinas asociadas no devuelven objetos sino listas de objetos. Esto da lugar a una compactación del AAA. Veáanse los códigos de `statements` y `idlist`:

```

sub statements() {
    my @s;

    @s = (statement());
    if ($lookahead eq ';'') {
        match(';');
        push @s, statements();
    }
    return @s;
}

sub idlist() {
    my @id;

    if ($lookahead eq 'ID') {
        @id = ($value); # no es un objeto
        match('ID');
        if ($lookahead eq ',') {
            match(',');
            push @id, idlist();
        }
    }
    else {
        Error::fatal('Se esperaba un identificador');
        @id = ('ERROR');
    }
    return @id;
}

```

4.9.4. Declaraciones

Los nodos del tipo declaration no existen propiamente, son nodos de la clase *INT* o de la clase *STRING*. La parte de la gramática árbol de la que hablamos es:

2	$\left \begin{array}{ll} \text{decls} & \rightarrow \text{list } \text{decl} \end{array} \right.$
4	$\left \begin{array}{ll} \text{decl} & \rightarrow \text{INT}(\text{idlist}) \end{array} \right.$
5	$\left \begin{array}{ll} \text{decl} & \rightarrow \text{STRING}(\text{idlist}) \end{array} \right.$
6	$\left \begin{array}{ll} \text{idlist} & \rightarrow \text{list } \text{SIMPLEID} \end{array} \right.$

Los nodos declaration son un hash con una clave TYPE la cual apunta a la estructura de datos/objeto describiendo el tipo. La otra clave del hash IDLIST apunta a una lista de identificadores. Los elementos de esta lista son simples identificadores (identificados en la gramática árbol anterior como *SIMPLEID* y no como objetos *ID*). La parte de la gramática implicada en las declaraciones es:

declaration \rightarrow INT idlist | STRING idlist
idlist \rightarrow ID ' ' idlist | ID

Así pues, el código construye un nodo de la clase INT o STRING según sea el caso.

```
sub declaration() {  
  my ($t, $class, @il);  
  
  if (($lookahead eq 'INT') or ($lookahead eq 'STRING')) {  
    $class = $lookahead;  
    $t = &type();  
    @il = &idlist();  
    return $class->new(TYPE => $t, IDLIST => \@il);  
  }  
  else {  
    Error::fatal('Se esperaba un tipo');  
  }  
}
```

Observe la llamada `$class->new(TYPE => $t, IDLIST => \@il)` en la cual la clase se usa a través de una referencia simbólica.

4.9.5. Práctica: Arbol de Análisis Abstracto

Complete la fase de análisis sintáctico para la gramática de Tutu extendida con sentencias de bloque (vea las reglas 1,2,3 y 11) construyendo el AAA según el lenguaje árbol especificado por una gramática árbol que extienda la dada en la definición 4.9.7. Genere pruebas, usando `make test` para comprobar el correcto funcionamiento de su analizador sobre las mismas. Utilize el módulo `Data::Dumper` para volcar las estructuras de datos resultantes.

1	$p \rightarrow b$
2	$b \rightarrow ds\ ss$
3	$b \rightarrow ss$
4	$ds \rightarrow d\ ';\ ' ds$
5	$ds \rightarrow d\ ';\ '$
6	$d \rightarrow INT\ il$
7	$d \rightarrow STRING\ il$
8	$ss \rightarrow s\ ';\ ' ss$
9	$ss \rightarrow s$
10	$s \rightarrow ID = e$
11	$s \rightarrow \{'\ b\ '\}$
12	$s \rightarrow P\ e$
13	$s \rightarrow \epsilon$
14	$e \rightarrow e1\ '+'\ t$
15	$e \rightarrow e1\ '\cdot'\ t$
16	$e \rightarrow t$
17	$t \rightarrow t1\ '\cdot'\ f$
18	$t \rightarrow t\ '/'\ f$
19	$t \rightarrow f$
20	$f \rightarrow '\(' e '\)'$
21	$f \rightarrow ID$
22	$f \rightarrow NUM$
23	$f \rightarrow STR$
24	$il \rightarrow ID\ '\ ' il$
25	$il \rightarrow ID$

4.10. Análisis Semántico

Hay quien dice que el análisis semántico es la determinación de aquellas propiedades que, siendo dependientes del contexto, pueden ser computadas estáticamente en tiempo de compilación para cualquier programa correcto. Entre estas propiedades están: la comprobación de que las variables son declaradas, la compatibilidad de tipos en las expresiones, el correcto uso de las llamadas a función así como el ámbito y visibilidad de las variables. La fase de análisis semántico puede verse como una fase de “adornado” o “etiquetado” del AAA, en la cual los atributos de los nodos del AAA son computados.

Aunque la veamos como una fase separada del análisis sintáctico, puede en numerosas ocasiones llevarse a cabo al mismo tiempo que se construye el árbol. Así lo hacemos en este ejemplo: incrustamos la acción semántica en la correspondiente rutina de análisis sintáctico. Así, en la rutina `term`, una vez que hemos obtenido los dos operandos, comprobamos que son de tipo numérico llamando (línea 8) a

`Semantic::Analysis::check_type_numeric_operator:`

Observe como aparece un nuevo atributo `TYPE` decorando el nodo creado (línea 9):

```

1 sub term() {
2   my ($t, $t2);
3
4   $t = factor;
5   if ($lookahead eq '*') {
6     match('*');
7     $t2 = term;
8     my $type = Semantic::Analysis::check_type_numeric_operator($t, $t2, '*');
9     $t = TIMES->new( LEFT => $t, RIGHT => $t2, TYPE => $type);
10  }
11  return $t;
12 }
```

En el manejo de errores de tipo, un tipo especial `$err_type` es usado para indicar un error de tipo:

```
sub check_type_numeric_operator {
    my ($op1, $op2, $operator) = @_;

    my $type = numeric_compatibility($op1, $op2, $operator);
    if ($type == $err_type) {
        Error::fatal("Operandos incompatibles para el operador $operator")
    }
    else {
        return $type;
    }
}
```

La subrutina `numeric_compatibility` comprueba que los dos operandos son de tipo numérico y devuelve el correspondiente tipo. Si ha ocurrido un error de tipo, intenta encontrar un tipo conveniente para el operando:

```
sub numeric_compatibility {
    my ($op1, $op2, $operator) = @_;

    if (($op1->TYPE == $op2->TYPE) and is_numeric($op1->TYPE)) {
        return $op1->TYPE; # correct
    }
    ... # código de recuperación de errores de tipo
}

sub is_numeric {
    my $type = shift;

    return ($type == $int_type); # añadir long, float, double, etc.
}
```

Es parte del análisis semántico la declaración de tipos:

```
sub declaration() {
    my ($t, $class, @il);

    if (($lookahead eq 'INT') or ($lookahead eq 'STRING')) {
        $class = $lookahead;
        $t = &type();
        @il = &idlist();
        &Semantic::Analysis::set_types($t, @il);
        &Address::Assignment::compute_address($t, @il);
        return $class->new(TYPE => $t, IDLIST => \@il);
    }
    else {
        Error::fatal('Se esperaba un tipo');
    }
}
```

Para ello se utiliza una tabla de símbolos que es un hash `%symbol_table` indexado en los identificadores del programa:

```

sub set_types {
  my $type = shift;
  my @vars = @_;

  foreach my $var (@vars) {
    if (!exists($symbol_table{$id})) { $symbol_table{$var}->{TYPE} = $type; }
    else { Error::error("$id declarado dos veces en el mismo ámbito"); }
  }
}

```

Cada vez que aparece una variable en el código, bien en un factor o en una asignación, comprobamos que ha sido declarada:

```

sub factor() {
  my ($e, $id, $str, $num);

  if ($lookahead eq 'NUM') { ... }
  elsif ($lookahead eq 'ID') {
    $id = $value;
    match('ID');
    my $type = Semantic::Analysis::check_declared($id);
    return ID->new( VAL => $id, TYPE => $type);
  }
  elsif ($lookahead eq 'STR') { ... }
  elsif ($lookahead eq '(') { ... }
  else { Error::fatal("Se esperaba (, NUM o ID"); }
}

```

La función `check_declared` devuelve el atributo `TYPE` de la correspondiente entrada en la tabla de símbolos.

```

sub check_declared {
  my $id = shift;

  if (!exists($symbol_table{$id})) {
    Error::error("$id no ha sido declarado!");
    # auto-declaración de la variable a err_type
    Semantic::Analysis::set_types($err_type, ($id));
  }
  return $symbol_table{$id}->{TYPE};
}

```

4.10.1. Práctica: Declaraciones Automáticas

Modifique la subrutina `check_declared` para que cuando una variable no haya sido declarada se declare “sobre la marcha”. ¿Puede utilizar información dependiente del contexto para decidir cual es la mejor forma de declararla?

4.10.2. Práctica: Análisis Semántico

Extienda el código de la práctica 4.9.5 para comprobar la compatibilidad de tipos.

1	$p \rightarrow b$
2	$b \rightarrow ds\ ss$
3	$b \rightarrow ss$
4	$ds \rightarrow d\ ';\ ' ds$
5	$ds \rightarrow d\ ';\ '$
6	$d \rightarrow INT\ il$
7	$d \rightarrow STRING\ il$
8	$ss \rightarrow s\ ';\ ' ss$
9	$ss \rightarrow s$
10	$s \rightarrow ID\ '=' e$
11	$s \rightarrow '\{ ' b ' \}'$
12	$s \rightarrow P\ e$
13	$s \rightarrow \epsilon$
14	$e \rightarrow el\ '+' t$
15	$e \rightarrow el\ '\cdot' t$
16	$e \rightarrow t$
17	$t \rightarrow t1\ '*' f$
18	$t \rightarrow t\ '/' f$
19	$t \rightarrow f$
20	$f \rightarrow '(' e ')'$
21	$f \rightarrow ID$
22	$f \rightarrow NUM$
23	$f \rightarrow STR$
24	$il \rightarrow ID\ ',' il$
25	$il \rightarrow ID$

En cuanto a las sentencias de bloque, se pretende que el ámbito y visibilidad de las variables sea como en el lenguaje C, esto es, las declaraciones mas internas con el mismo identificador ocultan las mas externas. Así:

```
int a;
a = 4;
{
    int a;
    a = 5;
    p a
}; /* el ; es necesario */
p a
```

Imprimiría 5 y 4. Para traducir esta sentencia es necesario *usar una lista/pila de referencias a tablas de símbolos. Cada sentencia compuesta o bloque tendrá su propia tabla de símbolos*. Los identificadores se búscan en la lista de referencias a tablas de símbolos, primero en la última tabla de símbolos insertada y sino se encuentra se busca en la penúltima insertada, etc.

Guarde como un atributo del identificador (*SYMTABLE*) la referencia a la tabla de símbolos a la que pertenece. Guarde como un atributo del nodo bloque (*BLOCK*) la referencia a la tabla de símbolos asociada.

4.11. Optimización Independiente de la Máquina

En esta fase se hace un análisis del árbol, sometiéndolo a transformaciones que aumenten la eficiencia del código final producido.

Ejemplos de tareas que se pueden llevar a cabo en esta fase son:

- Extracción del interior de un bucle de cálculos que son invariantes del bucle
- Plegado de constantes: computar las expresiones constantes en tiempo de compilación, no de ejecución

- Propagación de las constantes: si se sabe que una variable en un punto del programa tiene un valor constante `a = 4`, se puede sustituir su uso por el de la constante
- Eliminación de computaciones redundantes, cuando la misma expresión aparece repetidas veces con los mismos valores de las variables
- La eliminación de código “muerto”: código que se sabe que nunca podrá ser ejecutado

En nuestro primer ejemplo, reduciremos esta fase a realizar una tarea de plegado de las constantes. Primero lo haremos mediante la rutina

```
&Machine::Independent::Optimization::fold
```

En esta fase transformamos los AAA: si tenemos un árbol de la forma `OPERATION(left, right)`, esto es, su raíz es una operación, primero plegamos los subárboles `left` y `right`, y si se han transformado en constantes numéricas, entonces plegamos el nodo que pasa a ser numérico:

```
1 sub operator_fold { # Obsérvese el uso del aliasing
2
3   if ($_[0]->LEFT->is_operation) {
4     $_[0]->{LEFT}->fold;
5   }
6   if ($_[0]->RIGHT->is_operation) {
7     $_[0]->{RIGHT}->fold;
8   }
9   if (ref($_[0]->LEFT) eq "NUM" and ref($_[0]->RIGHT) eq "NUM") {
10    $_[0] = reduce_children($_[0]);
11  }
12 }
13
14 sub PLUS::fold {
15   operator_fold(@_);
16 }
17
18 sub TIMES::fold {
19   operator_fold(@_);
20 }
```

El plegado de las operaciones binarias se ha delegado en la subrutina `operator_fold`. En las líneas 3 y 6 se comprueba que se trata de un nodo de tipo operación. Si es así se procede a su plegado. Una vez plegados los dos subárboles hijo comprobamos en la línea 9 que los hijos actuales son de la clase `NUM`. Si es el caso, en la línea 10 cambiamos el nodo por el resultado de operar los dos hijos. Los nodos han sido extendidos con un método `is_operation` que determina si se trata de un nodo operación binaria o no. Para ello se han introducido nuevas clases de nodos: la clase `Node` está en la raíz de la jerarquía de herencia, las clases `Leaf` y `Binary` se usan para representar los nodos hoja y binarios y heredan de la anterior. Una clase informa a Perl que desea heredar de otra clase añadiendo el nombre de esa clase a la variable `@ISA` de su paquete. La herencia en Perl determina la forma de búsqueda de un método. Si el objeto no se puede encontrar en la clase, recursivamente y en orden primero-profundo se busca en las clases de las cuales esta hereda, esto es en las clases especificadas en el vector `@ISA`.

```
package Node;
```

```
sub is_operation {
  my $node = shift;
```

```

    return ref($node) =~ /^(TIMES)|(PLUS)$/;
}

package Leaf; # hoja del AAA
our @ISA = ("Node");
sub children {
    return ();
}

package Binary;
our @ISA = ("Node");
sub children {
    my $self = shift;

    return (LEFT => $self->{LEFT}, RIGHT => $self->{RIGHT});
}

```

Así pues, los objetos de la clase `Leaf` tienen acceso al método `is_operation`.

Ahora hacemos que las clases `PLUS` y `TIMES` hereden de la clase `BINARY`:

```

package PLUS;
our @ISA = ("Binary");

sub operator {
    my $self = shift;

    $_[0]+$_[1];
}

....

package TIMES;
our @ISA = ("Binary");

sub operator {
    my $self = shift;

    $_[0]*$_[1];
}

....

```

Obsérvese que en las líneas 4 y 7 del código del plegado de nodos de operación se ha accedido directamente al dato en vez de usar el método para modificar el atributo, saltándonos lo que la buena programación orientada a objetos indica. La forma en la que esta escrito hace que, por ejemplo, `$_[0]->{LEFT}` sea modificado. Recuérdesse que en Perl *los argumentos son alias de los parámetros*.

La subrutina `reduce_children` es la encargada de crear el nuevo nodo con el resultado de operar los hijos izquierdo y derecho:

```

1 sub reduce_children {
2     my ($node) = @_;
3
4     my $value = $node->operator($node->LEFT->VAL, $node->RIGHT->VAL);
5     NUM->new(VAL => $value, TYPE => $PL::Tutu::int_type);
6 }

```


En la línea 4 se usa el método `operator` asociado con un nodo operación.

Plegar una sentencia de impresión es plegar la expresión a imprimir:

```
sub PRINT::fold {  
    $_[0]->{EXPRESSION}->fold;  
}
```

Plegar una sentencia de asignación es plegar la parte derecha de la asignación:

```
sub ASSIGN::fold {  
    $_[0]->{RIGHT}->fold;  
}
```

de nuevo, hemos accedido a los campos en vez de usar los métodos.

Las restantes operaciones de plegado son triviales:

```
sub ID::fold { }
```

```
sub NUM::fold { }
```

```
sub STR::fold { }
```

Por último, para plegar todas las expresiones recorreremos la lista de sentencias del programa y las plegamos una a una.

```
sub fold {  
    my @statements = @{$tree->{STS}};  
    for my $s (@statements) {  
        $s->fold;  
    }  
}
```

4.11.1. Práctica: Plegado de las Constantes

Complete su proyecto de compilador de Tutu con la fase de plegado de las constantes siguiendo la metodología explicada en los párrafos previos. Mejore la jerarquía de clases con una clase abstracta `Operation` que represente a los nodos que se corresponden con operaciones binarias. Defina el método abstracto `operation` en dicha clase. Un *método abstracto* es uno que, mas que proveer un servicio representa un servicio o categoría. La idea es que al definir un clase base abstracta se indica un conjunto de métodos que deberían estar definidos en todas las clases que heredan de la clase base abstracta. Es como una declaración de interfaz que indica la necesidad de definir su funcionalidad en las clases descendientes, pero que no se define en la clase base. Un método abstracto debe producir una excepción con el mensaje de error adecuado si no se ha redefinido en la clase descendiente.

Para ello use la clave `abstract` del módulo `Class::MethodMaker`. Consulte la documentación del módulo `Class::MethodMaker`. Consulte [?] para saber más sobre clases abstractas.

4.12. Patrones Árbol y Transformaciones Árbol

En la fase de optimización presentada en la sección 4.11 transformabamos el programa en su representación intermedia, como un AAA decorado, para obtener otro AAA decorado.

Una transformación de un programa puede ser descrita como un conjunto de *reglas de transformación* o *esquema de traducción árbol* sobre el árbol abstracto que representa el programa.

Antes de seguir, es conveniente que repase los conceptos en la sección 4.9.1 sobre lenguajes y gramáticas árbol.

En su forma mas sencilla, estas reglas de transformación vienen definidas por ternas $(p, e, action)$, donde la primera componente de la terna p es un *patrón árbol* que dice que árboles deben ser seleccionados. La segunda componente e dice cómo debe transformarse el árbol que casa con el patrón p .

La acción *action* indica como deben computarse los atributos del árbol transformado a partir de los atributos del árbol que casa con el patrón p . Una forma de representar este esquema sería:

$$p \Longrightarrow e \{ \text{action} \}$$

Por ejemplo:

$$PLUS(NUM_1, NUM_2) \Longrightarrow NUM_3 \{ \$NUM_3\{VAL\} = \$NUM_1\{VAL\} + \$NUM_2\{VAL\} \}$$

cuyo significado es que dondequiera que haya un nódulo del AAA que case con el *patrón de entrada* $PLUS(NUM, NUM)$ deberá sustituirse el subárbol $PLUS(NUM, NUM)$ por el subárbol NUM . Al igual que en los esquemas de traducción, enumeramos las apariciones de los símbolos, para distinguirlos en la parte semántica. La acción indica como deben recomputarse los atributos para el nuevo árbol: El atributo **VAL** del árbol resultante es la suma de los atributos **VAL** de los operandos en el árbol que ha casado. La transformación se repite hasta que se produce la *normalización del árbol*.

Las reglas de “casamiento” de árboles pueden ser mas complejas, haciendo alusión a propiedades de los atributos, por ejemplo

$$ASSIGN(LEFTVALUE, x) \text{ and } \{ \text{notlive}(\$LEFTVALUE\{VAL\}) \} \Longrightarrow NIL$$

indica que se pueden eliminar aquellos árboles de tipo asignación en los cuáles la variable asociada con el nodo $LEFTVALUE$ no se usa posteriormente.

Otros ejemplos con variables S_1 y S_2 :

$$\begin{aligned} IFELSE(NUM, S_1, S_2) \text{ and } \{ \$NUM\{VAL\} \neq 0 \} &\Longrightarrow S_1 \\ IFELSE(NUM, S_1, S_2) \text{ and } \{ \$NUM\{VAL\} == 0 \} &\Longrightarrow S_2 \end{aligned}$$

Observe que en el patrón de entrada $ASSIGN(LEFTVALUE, x)$ aparece un “comodín”: la variable-árbol x , que hace que el árbol patrón $ASSIGN(LEFTVALUE, x)$ case con cualquier árbol de asignación, independientemente de la forma que tenga su subárbol derecho.

Las siguientes definiciones formalizan una aproximación simplificada al significado de los conceptos *patrones árbol* y *casamiento de árboles*.

Definición 4.12.1. Sea (Σ, ρ) un alfabeto con función de aridez y un conjunto (puede ser infinito) de variables $V = \{x_1, x_2, \dots\}$. Las variables tienen aridez cero:

$$\rho(x) = 0 \quad \forall x \in V.$$

Un elemento de $B(V \cup \Sigma)$ se denomina patrón sobre Σ .

Definición 4.12.2. Se dice que un patrón es un patrón lineal si ninguna variable se repite.

Definición 4.12.3. Se dice que un patrón es de tipo (x_1, \dots, x_k) si las variables que aparecen en el patrón leídas de izquierda a derecha en el árbol son x_1, \dots, x_k .

Ejemplo 4.12.1. Sea $\Sigma = \{A, CONS, NIL\}$ con $\rho(A) = \rho(NIL) = 0, \rho(CONS) = 2$ y sea $V = \{x\}$. Los siguientes árboles son ejemplos de patrones sobre Σ :

$$\{ x, CONS(A, x), CONS(A, CONS(x, NIL)), \dots \}$$

El patrón $CONS(x, CONS(x, NIL))$ es un ejemplo de patrón no lineal. La idea es que un patrón lineal como éste “fuerza” a que los árboles t que casen con el patrón deben tener iguales los dos correspondientes subárboles $t/1$ y $t/2$.¹ situados en las posiciones de las variables¹

Ejercicio 4.12.1. Dado la gramática árbol:

$$\begin{aligned} S &\rightarrow S_1(a, S, b) \\ S &\rightarrow S_2(NIL) \end{aligned}$$

¹Repase la notación de Dewey introducida en la definición 4.9.6

la cuál genera los árboles concretos para la gramática

$$S \rightarrow aSb \mid \epsilon$$

¿Es $S_1(a, X(NIL), b)$ un patrón árbol sobre el conjunto de variables $\{X, Y\}$? ¿Lo es $S_1(X, Y, a)$? ¿Es $S_1(X, Y, Y)$ un patrón árbol?

Ejemplo 4.12.2. Ejemplos de patrones para el AAA definido en el ejemplo 4.9.2 para el lenguaje Tutu son:

$$x, y, PLUS(x, y), ASSIGN(x, TIMES(y, ID)), PRINT(y) \dots$$

considerando el conjunto de variables $V = \{x, y\}$. El patrón $ASSIGN(x, TIMES(y, ID))$ es del tipo (x, y) .

Definición 4.12.4. Una sustitución es una aplicación θ que asigna variables a patrones $\theta : V \rightarrow B(V \cup \Sigma)$.

Tal función puede ser naturalmente extendida de las variables a los árboles: los nodos (hoja) etiquetados con dichas variables son sustituidos por los correspondientes subárboles.

$$\begin{aligned} \theta &: B(V \cup \Sigma) \rightarrow B(V \cup \Sigma) \\ t\theta &= \begin{cases} x\theta & \text{si } t = x \in V \\ a(t_1\theta, \dots, t_k\theta) & \text{si } t = a(t_1, \dots, t_k) \end{cases} \end{aligned}$$

Obsérvese que, al revés de lo que es costumbre, la aplicación de la sustitución θ al patrón se escribe por detrás: $t\theta$.

También se escribe $t\theta = t\{x_1/x_1\theta, \dots, x_k/x_k\theta\}$ si las variables que aparecen en t de izquierda a derecha son x_1, \dots, x_k .

Ejemplo 4.12.3. Si aplicamos la sustitución $\theta = \{x/A, y/CONS(A, NIL)\}$ al patrón $CONS(x, y)$ obtenemos el árbol $CONS(A, CONS(A, NIL))$. En efecto:

$$CONS(x, y)\theta = CONS(x\theta, y\theta) = CONS(A, CONS(A, NIL))$$

Ejemplo 4.12.4. Si aplicamos la sustitución $\theta = \{x/PLUS(NUM, x), y/TIMES(ID, NUM)\}$ al patrón $PLUS(x, y)$ obtenemos el árbol $PLUS(PLUS(NUM, x), TIMES(ID, NUM))$:

$$PLUS(x, y)\theta = PLUS(x\theta, y\theta) = PLUS(PLUS(NUM, x), TIMES(ID, NUM))$$

Definición 4.12.5. Se dice que un patrón $\tau \in B(V \cup \Sigma)$ con variables x_1, \dots, x_k casa con un árbol $t \in B(\Sigma)$ si existe una sustitución de τ que produce t , esto es, si existen $t_1, \dots, t_k \in B(\Sigma)$ tales que $t = \tau\{x_1/t_1, \dots, x_k/t_k\}$. También se dice que τ casa con la sustitución $\{x_1/t_1, \dots, x_k/t_k\}$.

Ejemplo 4.12.5. El patrón $\tau = CONS(x, NIL)$ casa con el árbol $t = CONS(CONS(A, NIL), NIL)$ y con el subárbol $t.1$. Las respectivas sustituciones son $t\{x/CONS(A, NIL)\}$ y $t.1\{x/A\}$.

$$\begin{aligned} t &= \tau\{x/CONS(A, NIL)\} \\ t.1 &= \tau\{x/A\} \end{aligned}$$

Ejercicio 4.12.2. Sea $\tau = PLUS(x, y)$ y $t = TIMES(PLUS(NUM, NUM), TIMES(ID, ID))$. Calcule los subárboles t' de t y las sustituciones $\{x/t_1, y/t_2\}$ que hacen que τ case con t' .

Por ejemplo es obvio que para el árbol raíz t/ϵ no existe sustitución posible:

$$t = TIMES(PLUS(NUM, NUM), TIMES(ID, ID)) = \tau\{x/t_1, y/t_2\} = PLUS(x, y)\{x/t_1, y/t_2\}$$

ya que un término con raíz $TIMES$ nunca podrá ser igual a un término con raíz $PLUS$.

El problema aquí es equivalente al de las expresiones regulares en el caso de los lenguajes lineales. En aquellos, los autómatas finitos nos proveen con un mecanismo para reconocer si una determinada cadena “casa” o no con la expresión regular. Existe un concepto análogo, el de *autómata árbol* que resuelve el problema del “casamiento” de patrones árbol. Al igual que el concepto de autómata permite la construcción de software para la búsqueda de cadenas y su posterior modificación, el concepto de autómata árbol permite la construcción de software para la búsqueda de los subárboles que casan con un patrón árbol dado.

Estamos ahora en condiciones de plantear una segunda aproximación al problema de la optimización independiente de la máquina utilizando una subrutina que busque por aquellos árboles que queremos optimizar (en el caso del plegado los árboles de tipo operación) y los transforme adecuadamente.

La función `match_and_transform_list` recibe una lista de árboles los cuales recorre sometiéndolos a las transformaciones especificadas. La llamada para producir el plegado sería:

```
Tree::Transform::match_and_transform_list(
  NODES => $tree->{STS}, # lista de sentencias
  PATTERN => sub {
    $_[0]->is_operation and $_[0]->LEFT->isa("NUM")
    and $_[0]->RIGHT->isa("NUM")
  },
  ACTION => sub {
    $_[0] = Machine::Independent::Optimization::reduce_children($_[0])
  }
);
```

Además de la lista de nodos le pasamos como argumentos una referencia a la subrutina encargada de reconocer los patrones árbol (clave `PATTERN`) y una referencia a la subrutina que describe la acción que se ejecutará (clave `ACTION`) sobre el árbol que ha casado. Ambas subrutinas asumen que el primer argumento que se les pasa es la referencia a la raíz del árbol que está siendo explorado.

Los métodos *isa*, *can* y *VERSION* son proporcionados por una clase especial denominada clase `UNIVERSAL`, de la cual implícitamente hereda toda clase. El método *isa* nos permite saber si una clase hereda de otra.

La subrutina `match_and_transform_list` recibe los argumentos y da valores por defecto a los mismos en el caso de que no hayan sido establecidos. Finalmente, llama a `match_and_transform` sobre cada uno de los nodos “sentencia” del programa.

```
sub match_and_transform_list {
  my %arg = @_;
  my @statements = @{$arg{NODES}} or
    Error::fatal("Internal error.match_and_transform_list ".
      "espera una lista anónima de nodos");
  local $pattern = ($arg{PATTERN} or sub { 1 });
  local @pattern_args = @{$arg{PATTERN_ARGS}} if defined $arg{PATTERN_ARGS};
  local $action = ($arg{ACTION} or sub { print ref($_),"\n" });
  local @action_args = @{$arg{ACTION_ARGS}} if defined $arg{ACTION_ARGS};

  for (@statements) {
    match_and_transform($_);
  }
}
```

La subrutina `match_and_transform` utiliza el método *can* para comprobar que el nodo actual dispone de un método para calcular la lista con los hijos del nodo. Una vez transformados los subárboles del nodo actual procede a comprobar que el nodo casa con el patrón y si es el caso le aplica la acción definida:

```

package Tree::Transform;

our $pattern;
our @pattern_args;
our $action;
our @action_args;
our @statements;

sub match_and_transform {
    my $node = $_[0] or Error::fatal("Error interno. match_and_transform necesita un nodo");
    Error::fatal("Error interno. El nodo de la clase",ref($node),
        " no dispone de método 'children'") unless $node->can("children");

    my %children = $node->children;

    for my $k (keys %children) {
        $node->{$k} = match_and_transform($children{$k});
    }

    if ($pattern->($node, @pattern_args)) {
        $action->($node, @action_args);
    }
    return $node;
}

```

Recordemos el esquema de herencia que presentamos en la sección anterior. Las clases Leaf y Binary proveen versiones del método children. Teníamos:

```

package Node;

sub is_operation {
    my $node = shift;

    return ref($node) =~ /^(TIMES)|(PLUS)$/;
}

package Leaf; # hoja del AAA
our @ISA = ("Node");
sub children {
    return ();
}

package Binary;
our @ISA = ("Node");
sub children {
    my $self = shift;

    return (LEFT => $self->{LEFT}, RIGHT => $self->{RIGHT});
}

```

Los objetos de la clase Leaf tienen acceso al método is_operation.

Las clases PLUS y TIMES heredan de la clase BINARY:

```

package PLUS;
our @ISA = ("Binary");

```

```

sub operator {
    my $self = shift;

    $_[0]+$_[1];
}

....

package TIMES;
our @ISA = ("Binary");

sub operator {
    my $self = shift;

    $_[0]*$_[1];
}

....

```

La subrutina `reduce_children` introducida en la sección 4.11 es la encargada de crear el nuevo nodo con el resultado de operar los hijos izquierdo y derecho:

```

1 sub reduce_children {
2     my ($node) = @_;
3
4     my $value = $node->operator($node->LEFT->VAL, $node->RIGHT->VAL);
5     NUM->new(VAL => $value, TYPE => $PL::Tutu::int_type);
6 }

```

En la línea 4 se usa el método `operator` asociado con un nodo operación.

4.12.1. Práctica: Casando y Transformando Árboles

Complete su proyecto para el compilador de Tutu completando las subrutinas `match_and_transform` tal y como se explicó en la sección 4.12.

Ademas del plegado de constantes use las nuevas subrutinas para aplicar simultáneamente las siguientes transformaciones algebraicas:

$$\begin{array}{lll}
 PLUS(NUM, x) & \wedge \{ \$NUM\{VAL\} == 0 \} & \implies x \\
 PLUS(x, NUM) & \wedge \{ \$NUM\{VAL\} == 0 \} & \implies x \\
 TIMES(NUM, x) & \wedge \{ \$NUM\{VAL\} == 1 \} & \implies x \\
 TIMES(x, NUM) & \wedge \{ \$NUM\{VAL\} == 1 \} & \implies x
 \end{array}$$

1. Dado un programa como

```
int a; a = a * 4 * 5;
```

¿Será plegado el `4 * 5`? Sin embargo si que se pliega si el programa es de la forma:

```
int a; a = a * (4 * 5);
```

No intente en esta práctica que programas como el primero o como `4*a*5*b` sean plegados. Para lograrlo sería necesario introducir transformaciones adicionales y esto no se requiere en esta práctica.

2. ¿Existe un orden óptimo en el que ejecutar las transformaciones?
3. Ponga un ejemplo en el que sea beneficioso ejecutar el plegado primero.
4. Ponga otro ejemplo en el que sea beneficioso ejecutar el plegado después.
5. ¿Es necesario aplicar las transformaciones reiteradamente?
6. ¿Cuál es la condición de parada?
7. Como es habitual la pregunta 6 tiene una respuesta TIMTOWTDI: una posibilidad la da el módulo `Data::Compare` el cual puede obtenerse desde CPAN y que permite comparar estructuras de datos, pero existe una solución mas sencilla. ¿Cuál?

4.13. Asignación de Direcciones

Esta suele ser considerada la primera de las fases de síntesis. Las anteriores lo eran de análisis. La fase de análisis es una transformación *texto fuente* \rightarrow *arbol*, mientras que la fase síntesis es una transformación inversa *arbol* \rightarrow *texto objeto* que produce una “linealización” del árbol. En general, se ubican en la fase de síntesis todas las tareas que dependan del lenguaje objeto.

La asignación de direcciones depende de la máquina objetivo en cuanto conlleva consideraciones sobre la longitud de palabra, las unidades de memoria direccionables, la compactación de objetos pequeños (por ejemplo, valores lógicos), el alineamiento a fronteras de palabra, etc.

En nuestro caso debemos distinguir entre las cadenas y las variables enteras.

Las constantes literales (como “hola”) se almacenan concatenadas en orden de aparición textual. Una variable de tipo cadena ocupa dos palabras, una dando su dirección y otra dando su longitud.

```
sub factor() {
    my ($e, $id, $str, $num);

    if ($lookahead eq 'NUM') { ... }
    elsif ($lookahead eq 'ID') { ... }
    elsif ($lookahead eq 'STR') {
        $str = $value;
        my ($offset, $length) = Address::Assignment::str($str);
        match('STR');
        return STR->new(OFFSET => $offset, LENGTH => $length, TYPE => $string_type);
    }
    elsif ($lookahead eq '(') { ... }
    else { Error::fatal("Se esperaba (, NUM o ID"); }
}
```

El código de la subrutina `Address::Assignment::str` es:

```
sub str {
    my $str = shift;
    my $len = length($str);
    my $offset = length($data);
    $data .= $str;
    return ($offset, $len);
}
```

Una posible mejora es la que se describe en el punto 2 de la práctica 4.13.1.

Hemos supuesto que las variables enteras y las referencias ocupan una palabra. Cada vez que una variable es declarada, se le computa su dirección relativa:

```

# declaration -> type idlist
# type          -> INT | STRING
sub declaration() {
    my ($t, $decl, @il);

    if (($lookahead eq 'INT') or ($lookahead eq 'STRING')) {
        $t = &type(); @il = &idlist();
        &Semantic::Analysis::set_types($t, @il);
        &Address::Assignment::compute_address($t, @il);
        $decl = [$t, \@il];
        return bless $decl, 'DECL';
    }
    else { Error::fatal('Se esperaba un tipo'); }
}

```

Se usa una variable `$global_address` para llevar la cuenta de la última dirección utilizada y se introduce un atributo `ADDRESS` en la tabla de símbolos:

```

sub compute_address {
    my $type = shift;
    my @id_list = @_;

    for my $i (@id_list) {
        $symbol_table{$i}->{ADDRESS} = $global_address;
        $global_address += $type->LENGTH;
    }
}

```

Por último situamos todas las cadenas después de las variables del programa fuente:

```

...
##### En compile, despues de haber calculado las direcciones
Tree::Transform::match_and_transform_list(
    NODES => $tree->{STS},
    PATTERN => sub {
        $_[0]->isa('STR')
    },
    ACTION => sub { $_[0]->{OFFSET} += $global_address; }
);

```

Esta aproximación es bastante simplista en cuanto que usa una palabra por carácter. El punto 3 de la práctica 4.13.1 propone una mejora.

4.13.1. Práctica: Cálculo de las Direcciones

Modifique el cálculo de las direcciones para la gramática de Tutu extendida con sentencias compuestas que fué presentada en la sección 4.10.2.

1. Resuelva el problema de calcular las direcciones de las variables declaradas en los bloques. La memoria ocupada por dichas variables se libera al terminar el bloque. Por tanto la variable `$global_address` disminuye al final de cada bloque.
2. En el código de la subrutina `Address::Assignment::str`


```

sub str {
    my $str = shift;
    my $len = length($str);
    my $offset = length($data);
    $data .= $str;
    return ($offset, $len);
}

```

no se comprueba si la cadena `$str` ya está en `$data`. Es natural que si ya está no la insertemos de nuevo. Introduzca esa mejora. Para ello use la función `pos`, la cual devuelve el desplazamiento en `$data` donde quedó la última búsqueda `$data =~ m/regex/g`. El siguiente ejemplo con el depurador le muestra como trabaja. Escriba `perldoc -f pos` para obtener mas información sobre la función `pos`. Vea un ejemplo de uso:

```

DB<1> $str = "world"
#          012345678901234567890123456789012
DB<2> $data = "hello worldjust another statement"
DB<3> $res = $data =~ m{$str}g
DB<4> x pos($data) # la siguiente busqueda comienza en la 11
0 11
DB<5> $s = "hello"
DB<6> $res = $data =~ m{$s}g
DB<7> x pos($data) # No hay "hello" despues de world en $data
0 undef
DB<8> $res = $data =~ m{$str}g # Repetimos ...
DB<9> x pos($data)
0 11
DB<10> pos($data) = 0 # Podemos reiniciar pos!
DB<11> $res = $data =~ m{$s}g
DB<12> x pos($data) # Ahora si se encuentra "hello"
0 5

```

3. Empaquete las cadenas asumiendo que cada carácter ocupa un octeto o byte. Una posible solución al problema de alineamiento que se produce es rellenar con ceros los bytes que faltan hasta completar la última palabra ocupada por la cadena. Por ejemplo, si la longitud de palabra de la máquina objeto es de 4 bytes, la palabra *procesador* se cargaría así:

```

    0   |   1   |   2
0 1 2 3 | 4 5 6 7 | 8 9 10 11
p r o c | e s a d | o r \0 \0

```

Existen dos posibles formas de hacer esta empaquetado. Por ejemplo, si tenemos las cadenas *lucas* y *pedro* las podemos introducir en la cadena `$data` así:

```

    0   |   1   |   2   |   3
0 1 2 3 | 4 5 6 7 | 8 9 10 11 | 12 13 14 15
l u c a | s \0 \0 \0 | p e d r | o \0 \0 \0

```

o bien:

```

    0   |   1   |   2
0 1 2 3 | 4 5 6 7 | 8 9 10 11
l u c a | s p e d | r o \0 \0

```

La segunda forma aunque compacta más tiene la desventaja de hacer luego mas lento el código para el direccionamiento de las cadenas, ya que no empiezan en frontera de palabra. Utilice el primer modo.

Observe que este empaquetado introduce un efecto en lo que se considera un éxito en la búsqueda de una cadena `$str` en `$data`. ¿Que ocurre si `$str` aparece en `$data` como subcadena de una cadena previamente empaquetada ocupando una posición que no es frontera de palabra?

Cuando rellene con `\0` la cadena use el operador `x` (letra equis) para multiplicar número por cadena:

```
DB<1> $x = "hola"x4
DB<2> p $x
holaholaholahola
```

esto le evitará escribir un bucle.

4. Mantenga atributos distintos en el nodo `STR` para el desplazamiento y longitud de `$str` en `$data` (en caracteres) y su dirección y tamaño final en memoria (en palabras).

4.14. Generación de Código: Máquina Pila

El generador de código emite el código para el lenguaje objeto, utilizando las direcciones calculadas en la fase anterior. Hemos simplificado esta fase considerando una máquina orientada a pila: una memoria de pila en la que se computan las expresiones, un segmento de datos en el que se guardan las variables y un segmento para guardar las cadenas. La estrategia utilizada es similar a la que vimos en el plegado de las constantes. Definimos un método para la traducción de cada clase de nodo del AAA. Entonces, para traducir el programa basta con llamar al correspondiente método:

```
$tree = Syntax::Analysis::parser;
... # otras fases
#####code generation
local $target = ""; # target code
$tree->translate;
...
```

en la cadena `$target` dejamos el código emitido. Cada método visita los hijos, traduciéndolos y añadiendo el código que fuera necesario para la traducción del nodo.

```
sub PROGRAM::translate {
    my $tree = shift;

    $target .= "DATA ". $data."\n" if $data;
    $tree->STS->translate;
}
```

Traducir la lista de sentencias es concatenar la traducción de cada una de las sentencias en la lista:

```
sub STATEMENTS::translate {
    my $statements = shift;
    my @statements = @{$statements};

    for my $s (@statements) {
        $s->translate;
    }
}
```

Si suponemos que disponemos en el ensamblador de nuestra máquina objeto de instrucciones PRINT_INT y PRINT_STR que imprimen el contenido de la expresión que esta en la cima de la pila, la traducción será:

```
sub PRINT::translate {
    my $self = shift;

    $self->EXPRESSION->translate;
    if ($self->EXPRESSION->TYPE == $int_type) { emit "PRINT_INT\n"; }
    else {emit "PRINT_STR\n"; }
}
```

Asi, si la sentencia era P c, donde c es del tipo cadena, se debería eventualmente llamar al método de traducción del identificador:

```
sub ID::translate {
    my $self = shift;

    my $id = $self->VAL;
    my $type = Semantic::Analysis::get_type($id);
    if ($type == $int_type) {
        emit "LOAD ".$symbol_table{$id}->{ADDRESS}."\n";
    }
    else {
        emit "LOAD_STRING ".$symbol_table{$id}->{ADDRESS}."\n";
    }
}
```

la función emit simplemente concatena el código producido a la salida:

```
sub emit { $target .= shift; }
```

Para la traducción de una sentencia de asignación supondremos de la existencia de instrucciones STORE_INT y STORE_STRING. La instrucción STORE_STRING asume que en la cima de la pila están la dirección y la cadena a almacenar.

```
sub ASSIGN::translate {
    my $self = shift;

    $self->RIGHT->translate;
    my $id = $self->LEFT;
    $id->translate;
    my $type = Semantic::Analysis::get_type($id->VAL);

    if ($type == $int_type) {
        emit "STORE_INT\n";
    }
    else {
        emit "STORE_STRING\n";
    }
}
```

Si se está traduciendo una sentencia como a = "hola", se acabará llamando al método de traducción asociado con la clase STR el cual actúa empujando la dirección y la longitud de la cadena en la pila:

```

sub STR::translate {
    my $self = shift;

    emit "PUSHSTR ".$self->OFFSET." ".$self->LENGTH."\n";
}

```

Así la traducción de este fuente:

```

$ cat test06.tutu
string a;
a = "hola";
p a

```

es:

```

$ ./main.pl test06.tutu test06.ok
$ cat test06.ok
DATA 2, hola
PUSHSTR 2 4
PUSHADDR 0
STORE_STRING
LOAD_STRING 0
PRINT_STR

```

El resto de los métodos de traducción es similar:

```

sub LEFTVALUE::translate {
    my $id = shift ->VAL;

    emit "PUSHADDR ".$symbol_table{$id}->{ADDRESS}."\n";
}

```

```

sub NUM::translate {
    my $self = shift;

    emit "PUSH ".$self->VAL."\n";
}

```

```

sub PLUS::translate {
    my $self = shift;

    $self->LEFT->translate;
    $self->RIGHT->translate;
    emit "PLUS\n";
}

```

```

sub TIMES::translate {
    my $self = shift;

    $self->LEFT->translate;
    $self->RIGHT->translate;
    emit "MULT\n";
}

```

Veamos un ejemplo. Dado el código fuente:

```
$ cat test02.tutu
int a,b; string c;
a = 2+3; b = 3*4; c = "hola"; p c;
c = "mundo"; p c; p 9+2; p a+1; p b+1
```

el código generado es:

```
$ ./main.pl test02.tutu test02.ok
$ cat test02.ok
DATA 4, holamundo
PUSH 5
PUSHADDR 0
STORE_INT
PUSH 12
PUSHADDR 1
STORE_INT
PUSHSTR 4 4
PUSHADDR 2
STORE_STRING
LOAD_STRING 2
PRINT_STR
PUSHSTR 8 5
PUSHADDR 2
STORE_STRING
LOAD_STRING 2
PRINT_STR
PUSH 11
PRINT_INT
LOAD 0
INC
PRINT_INT
LOAD 1
INC
PRINT_INT
```

4.15. Generación de Código: Máquina Basada en Registros

La máquina orientada a pila para la que generamos código en la sección 4.14 es un ejemplo de la clase de máquinas que es usada por la mayoría de los lenguajes interpretados: Perl, Python; java, etc.

En esta sección introduciremos una máquina basada en registros. Suponemos que la máquina tiene k registros $R_0 \dots R_{k-1}$. Las instrucciones toman dos argumentos, dejando el resultado en el primer argumento. Son las siguientes:

LOADM R_i , $[a]$	$R_i = M_a$
LOADC R_i , c	$R_i = c$
STORE $[a]$, R_i	$M_a = R_i$
ADDR R_i , R_j	$R_i += R_j$
ADDM R_i , $[a]$	$R_i += M_a$
ADDC R_i , c	$R_i += c$
...	...

El problema es generar el código con el menor número de instrucciones posible, teniendo en cuenta la limitación existente de registros.

Supongamos que queremos traducir un subárbol $OP(t_1, t_2)$ y que la traducción del subárbol t_1 requiere r_1 registros y que la traducción de t_2 requiere r_2 registros, con $r_1 < r_2 \leq k$. Si realizamos

primero la evaluación de t_1 , debemos dejar el resultado en un registro que no podrá ser utilizado en la evaluación de t_2 . Si $r_2 = k$, la evaluación de t_2 podría dar lugar a la necesidad de recurrir a almacenamiento temporal. Esta situación no se da si evaluamos primero t_2 . En tal caso, dado que hay un registro en el que se guarda el resultado de t_2 , quedan libres al menos $r_2 - 1$ registros. Como $r_2 - 1 \geq r_1$ se sigue que tenemos suficientes registros para traducir t_1 . Como regla general es mejor evaluar primero el subárbol que mayores requerimientos de registros tiene.

La siguiente cuestión es como calcular los requerimientos en registros de una expresión dada. No consideraremos en esta fase límites en el número de registros disponibles. Obsérvese que si los requerimientos para los subárboles son distintos, $r_1 \neq r_2$ la traducción puede realizarse usando el máximo de ambos $\max\{r_1, r_2\}$ siguiendo la estrategia de traducir primero el que mayores requerimientos tenga. Si son iguales entonces se necesitan $r_1 + 1$ registros ya que es necesario un registro para guardar el resultado de la primera traducción.

Nótese que, como el juego de instrucciones para un operando puede tener como segundo argumento una dirección de memoria, los “segundos operandos” no necesitan registro. Por ejemplo, el árbol $PLUS(a, b)$ se traduce por

```
LOADM R0, a
PLUSM R0, b
```

Así b no requiere registro, mientras que a si lo requiere. Por tanto, las hojas izquierdas requieren de registro mientras que las hojas derechas no.

Si t es un nodo de la forma $OP(t_1, t_2)$ el número de registros r_t requeridos por t viene dado por la fórmula:

$$r_t = \begin{cases} \max\{r_1, r_2\} & \text{si } r_1 \neq r_2 \\ r_1 + 1 & \text{si } r_1 = r_2 \end{cases}$$

Dotaremos a cada nodo del AST de un método `required_registers` que computa la demanda en registros de dicho nodo. Lo que haremos es introducir en la clase `Operation` de la cual heredan las operaciones binarias el correspondiente método `required_registers`:

```
package Operation;
our @ISA = ("Binary");

sub required_registers {
    my $self = shift;

    my $rl = $self->LEFT->required_registers('LEFT');
    my $rr = $self->RIGHT->required_registers('RIGHT');
    $self->{REQ_REG} = ($rl == $rr)? $rl+1: Aux::max($rl, $rr);
    return $self->REQ_REG;
}
```

El segundo argumento que recibe `required_registers` es su posición (izquierda o derecha) entre los hijos de su padre. dicha información no es usada en los nodos binarios. Su necesidad queda clara cuando se considera el cómputo del número de registros requeridos por las hojas.

El cómputo en las hojas corre a cargo del correspondiente método en la clase `Value`. Los nodos de tipo número (clase `NUM`), cadena (clase `STR`) y variable (clase `ID`) heredan de la clase `Value`.

```
package Value;
our @ISA = ("Leaf");

sub required_registers {
    my $self = shift;
```

```

my $position = shift;

$self->{REQ_REG} = ($position eq 'LEFT') ? 1 : 0;
return $self->REQ_REG;
}

```

El atributo REQ_REG se computa para cada una de las sentencias del programa:

```

package STATEMENTS;

sub required_registers {
    my $self = shift;
    my @sts = @{$self};

    for (@sts) {
        $_->required_registers;
    }
}

```

Por supuesto los nodos ASSIGN y PRINT poseen sus propios métodos `required_registers`.

Una vez computados los requerimientos en registros de cada nodo, la generación de código para un nodo gestiona la asignación de registros usando una cola en la que se guardan los registros disponibles. Se siguen básicamente dos reglas para la traducción de un nodo **Operation**:

1. Realizar primero la traducción del hijo con mayores requerimientos y luego el otro
2. El resultado queda siempre en el registro que ocupa la primera posición en la cola

Hay cuatro casos a considerar: el primero es que el operando derecho sea una hoja. La generación de código para este caso es:

```

package Operation;
our @ISA = ("Binary");
...

sub gen_code {
    my $self = shift;

    if ($self->RIGHT->isa('Leaf')) {
        my $right = $self->RIGHT;
        my $a = $right->VAL;
        my $rightoperand = $right->gen_operand; # valor o dirección
        my $key = $right->key;                  # M, C, etc.
        $self->LEFT->gen_code;
        Aux::emit($self->nemonic."$key $RSTACK[0], $rightoperand # $a\n");
    }
    ...
}

```

La generación del nemónico se basa en tres métodos:

- El método `nemonic` devuelve el nemónico asociado con el nodo. Por ejemplo, para la clase **TIMES** el código es:

```

sub nemonic {
    return "MULT";
}

```

- El método `key` devuelve el sufijo que hay que añadir para completar el nemónico, en términos de como sea el operando: C para los números, M para los identificadores, etc.
- El método `gen_operand` genera el operando. Así para las clases número e identificador su código es:

<pre>package NUM; ... sub gen_operand { my \$self = shift; return \$self->VAL; }</pre>	<pre>package ID; ... sub gen_operand { my \$self = shift; return \$symbol_table{\$self->VAL}->{ADDRESS}, }</pre>
--	---

El resto del código distingue tres casos, según sean r_1 , r_2 y el número de registros disponibles. Los dos primeros casos desglosan la posibilidad de que uno de los dos subárboles pueda realizarse con el número de registros disponible ($\min\{r_1, r_2\} < k$). El tercer caso corresponde a que se necesiten temporales: $\min\{r_1, r_2\} \geq k$.

```
1  ...
2  if ($self->RIGHT->isa('Leaf')) { ... }
3  else { # Hijo derecho no es una hoja
4      my ($t1, $t2) = ($self->LEFT, $self->RIGHT);
5      my ($r1, $r2) = ($t1->REQ_REG, $t2->REQ_REG);
6
7      if ($r1 < Aux::min($r2, $NUM_REG)) {
8          $t2->gen_code;
9          my $R = shift @RSTACK;
10         $t1->gen_code;
11         Aux::emit($self->nemonic."R $RSTACK[0], $R\n");
12         push @RSTACK, $R;
13     }
14     ...
15 }
```

En este caso debemos realizar primero la traducción del hijo derecho. Salvando su resultado en `$R`. El registro es retirado de la cola y traducimos el lado izquierdo. El resultado ha quedado en el primer registro de la cola. Emitimos la operación, añadiendo el sufijo R, ya que se trata de una operación entre registros y posteriormente devolvemos el registro a la cola.

Ejercicio 4.15.1. *Responda a las siguientes preguntas:*

1. Si en el código anterior sustituimos la línea 12

`push @RSTACK, $R`

por

`unshift @RSTACK, $R`

¿Seguiría funcionando el código?

2. ¿Podemos asegurar en este subcaso que el código generado para \$t2 (línea 8) se ha realizado íntegramente en los registros?

Los otros dos casos tienen similar tratamiento:

```
if ($self->RIGHT->isa('Leaf')) { ... }
else { ...
  if ($r1 < Aux::min($r2, $NUM_REG)) { ... }
  elsif (($r1 >= $r2) and ($r2 < $NUM_REG)) {
    $t1->gen_code;
    my $R = shift @RSTACK;
    $t2->gen_code;
    Aux::emit($self->nemonic."R $R, $RSTACK[0]\n");
    unshift @RSTACK, $R;
  }
  elsif (($r1 >= $NUM_REG) and ($r2 >= $NUM_REG)) {
    $t2->gen_code;
    Aux::emit("STORE $T, $RSTACK[0]\n");
    $T++;
    $t1->gen_code;
    $T--;
    Aux::emit($self->nemonic."M $RSTACK[0], $T\n");
  }
}
}
```

Antes de comenzar a generar el código, la variable \$T debe ser inicializada a un valor apropiado, de manera que se usen direcciones no ocupadas por los datos. Por ejemplo:

```
local $T = $final_global_address+length($data);
```

El método `gen_code` sólo debería ser llamado sobre una hoja si se trata de una hoja izquierda (en cuyo caso el número de registros requeridos es uno):

```
package Value;
our @ISA = ("Leaf");
...

sub gen_code {
  my $self = shift;
  my $a = $self->VAL;

  if ($self->REQ_REG == 1) {
    if (ref($self) eq "NUM") { Aux::emit("LOADC $RSTACK[0], $a\n"); }
    else {
      my $address = $symbol_table{$a}->{ADDRESS};
      Aux::emit("LOADM $RSTACK[0], $address # $a\n");
    }
  }
  else {
    croak("gen_code visita hoja izquierda con REQ_REG = ".$self->REQ_REG);
  }
}
```

La pila de registros es inicializada al número de registros disponibles:

```
use constant LAST_REG => 1;
our @RSTACK = map "R$_", 0..LAST_REG; # Registros disponibles
```

Ejercicio 4.15.2. *Responda a las siguientes preguntas:*

- ¿Cuáles son los requerimientos de registros para un nodo de la clase *ASSIGN*?
 - ¿Cuáles son los requerimientos de registros para un nodo de la clase *PRINT*?
 - ¿Se puede lograr la funcionalidad proveída por el método `required_registers` usando `match_and_transform`?
- ¿Sería necesario introducir modificaciones en `match_and_transform`? Si es así, ¿Cuáles?.*

4.15.1. Práctica: Generación de Código

1. Complete la generación de código para la máquina basada en registros. Recuerde que debe escribir el método `required_registers` para las diferentes clases *Value*, *Operation*, *ASSIGN*, *PRINT*, *STATEMENTS*, etc. Así mismo deberá escribir el método `gen_code` para las diversas clases: *Value*, *Operation*, *ASSIGN*, *PRINT*, *STATEMENTS*, etc. Recuerde que los temporales usados durante la generación de código deben ubicarse en una zona que no esté en uso.
2. En la sección anterior no se consideraba la generación de código para las cadenas. Amplíe y/o modifique el juego de instrucciones como considere conveniente. El siguiente ejemplo de traducción sugiere como puede ser la generación de código para las cadenas:

Fuente	Objeto
<pre>string a,b; a = "hola"; b = a; p b</pre>	<pre>1 LSTRG R0, 4, 4 2 STORES 0, R0 # a 3 LOADS R0, 0 # a 4 STORES 2, R0 # b 5 LOADS R0, 2 # b 6 PRNTS R0</pre>

Asuma que los registros pueden contener dos direcciones de memoria (línea 1). La instrucción `LSTRG R0, a, b` carga las constantes (direcciones) `a` y `b` en el registro. La constante `"hola"` ocupa en la posición final en la que se colocan los contenidos de `$data` un desplazamiento de 4 y ocupa 4 palabras. Las instrucción `LOADS R, a` carga las dos palabras en las direcciones `a` y `a+1` en el registro `R`. La instrucción `STORES a, R` se encarga de que las dos palabras en la dirección `a` queden referenciando una cadena igual a la apuntada por el registro `R`. La instrucción `PRNTS` imprime la cadena apuntada por el registro. En una situación mas realista instrucciones como `STORES a, R` y `PRNTS` probablemente serían llamadas a funciones/servicios del sistema o de la librería para soporte en tiempo de ejecución asociada al lenguaje.

3. Se puede mejorar el código generado si hacemos uso de las propiedades algebraicas de los operadores. Por ejemplo, cuando se tiene un operador conmutativo que ha sido asociado a derechas, como ocurre en este programa fuente:

```
$ cat test18.tutu
int a,b,c;

a = a + (b + c)
```

El código producido por el compilador es:

```
LOADM R0, 0 # a
LOADM R1, 1 # b
PLUSM R1, 2 # c
PLUSR R0, R1
STORE 0, R0 # a
```

En este caso, la expresión $a + (b + c)$ corresponde a un árbol que casa con el patrón árbol

$$PLUS(ID, t) \text{ and } \{r_t \geq 1\}$$

Donde r_t es el número de registros requeridos por t . En tales casos es posible sacar ventaja de la conmutatividad de la suma y transformar el árbol

$$PLUS(ID, t) \text{ and } \{r_t \geq 1\} \implies PLUS(t, ID)$$

Observe que mientras el primer árbol requiere $\max\{2, r_t\}$ registros, el segundo requiere r_t registros, que en general es menor. Esta transformación invierte la traducción:

```
traduce(t)
ADDM $RSTACK[0], dirección de ID
```

que daría lugar a:

```
LOADM R0, 1 # b
PLUSM R0, 2 # c
PLUSM R0, 0 # a
STORE 0, R0 # a
```

la cual usa una instrucción y un registro menos.

Usando `match_and_transform` modifique el generador de código para que, después de la fase de cálculo del número de registros requeridos, aplique esta transformación sobre los nodos conmutativos cuyo hijo izquierdo sea un identificador y su hijo derecho requiera al menos un registro.

4.16. Optimización de Código

Aunque en esta fase se incluyen toda clase de optimizaciones, es aquí donde se hacen las optimizaciones de código dependientes del sistema objeto. Normalmente se recorre el código generado buscando secuencias de instrucciones que se puedan sustituir por otras cuya ejecución sea mas eficiente. El nombre *Peephole optimization* hace alusión a esta especie de “ventana de visión” que se desplaza sobre el código. En nuestro caso, supongamos que disponemos de una instrucción INC que permite incrementar eficientemente una expresión. Recorremos el código buscando por un patrón “sumar 1z lo reemplazamos adecuadamente.

```
package Peephole::Optimization;

sub transform {
    $target = shift;
    $target =~ s/PUSH 1\nPLUS/INC/g;
}
```

Otro ejemplo de optimización peephole consiste en reemplazar las operaciones flotantes de división por una constante por la multiplicación por la inversa de la misma (aquí se pueden introducir diferencias en el resultado, debido a la inexactitud de las operaciones en punto flotante y a que, si se trata de un *compilador cruzado* la aritmética flotante en la máquina en la que se ejecuta el compilador puede ser diferente de la de la máquina que ejecutará el código objeto).

4.16.1. Práctica: Optimización Peephole

1. Optimice el código generado para la máquina de registros sustituyendo las operaciones de multiplicación y división enteras por una constante que sea potencia de dos (de la forma 2^n) por operaciones de desplazamiento. Repase el capítulo de expresiones regulares. Es posible que aquí quiera emplear una sustitución en la que la cadena de reemplazo sea evaluada sobre la marcha. Si es así, repase la sección 3.1.6. La siguiente sesión con el depurador pretende ilustrar la idea:

```
$ perl -de 0
DB<1> $a = "MUL R2, 16"
DB<2> $a =~ s/MUL R(\d), (\d+)/($2 == 16)?"SHL R$1, 4":$&/e
DB<3> p $a
SHL R2, 4
DB<5> $a = "MUL R2, 7"
DB<6> $a =~ s/MUL R(\d), (\d+)/($2 == 16)?"SHL R$1, 4":$&/e
DB<7> p $a
MUL R2, 7
DB<8>
```

2. El plegado de constantes realizado durante la optimización de código independiente de la máquina (véase la sección 4.11) es parcial. Si los árboles para el producto se hunden a izquierdas, una expresión como $a = a * 2 * 3$ no será plegada, ya que produce un árbol de la forma

$$t = TIMES(TIMES(a, 2), 3)$$

Dado que el algoritmo no puede plegar $t/1$ tampoco plegará t . Busque en el código objeto secuencias de multiplicaciones por constantes y abrévelas en una. Haga lo mismo para las restantes operaciones.

3. Dado el siguiente programa:

```
$ cat test14.tutu
int a,b; a = 2; b = a*a+1
```

El código producido por el traductor para la máquina de registros es:

```
1 LOADC R0, 2
2 STORE 0, R0 # a
3 LOADM R0, 0 # a
4 MULTM R0, 0 # a
5 PLUSC R0, 1 # 1
6 STORE 1, R0 # b
```

Se ve que la instrucción de carga `LOADM R0, 0` de la línea 3 es innecesaria por cuanto el contenido de la variable `a` ya está en el registro `R0`, ya que fué cargada en el registro en la línea 2. Nótese que esta hipótesis no es necesariamente cierta si la línea 3 fuera el objetivo de un salto desde otro punto del programa. Esta condición se cumple cuando nos movemos dentro de un *bloque básico*: una secuencia de instrucciones que no contiene instrucciones de salto ni es el objetivo de instrucciones de salto, con la excepción de las instrucciones inicial y final. Mejore el código generado intentando detectar patrones de este tipo, eliminando la operación de carga correspondiente.

Capítulo 5

Construcción de Analizadores Léxicos

Habitualmente el término “análisis léxico” se refiere al tratamiento de la entrada que produce como salida la lista de *tokens*. Un *token* hace alusión a las unidades mas simples que tiene significado. Habitualmente un *token* o lexema queda descrito por una expresión regular. Léxico viene del griego *lexis*, que significa “palabra”. Perl es, sobra decirlo, una herramienta eficaz para encontrar en que lugar de la cadena se produce el emparejamiento. Sin embargo, en el análisis léxico, el problema es encontrar la subcadena a partir de la última posición “casada” que casa con una de las expresiones regulares que definen los lexemas del lenguaje dado.

5.1. Encontrando los terminales mediante sustitución

En esta sección construiremos una clase que provee una familia sencilla de analizadores léxicos. El constructor de la clase recibe como entrada la familia de parejas (expresión-regular, identificador-de-terminal) y devuelve como resultado una referencia a una subrutina, dinámicamente creada, que hace el análisis léxico. Esto es, se devuelve una subrutina que cuando se la llama con una cadena de entrada que se le pasa como parámetro devuelve la siguiente pareja (cadena, identificador-de-terminal).

Por ejemplo, después de la llamada:

```
my $lex = Lexer->new('\s*', #espacios
                    '\d+(\.\d+)?'=>'NUMBER',
                    '#.*'=>'COMMENT',
                    '"[^"]*"'=>>'STRING',
                    '\$[a-zA-Z_]\w*'=>'VAR',
                    '[a-zA-Z_]\w*'=>'ID',
                    '.'=>'OTHER');
```

la variable `$lex` contendrá una referencia a una subrutina que se creará dinámicamente y cuyo código es similar al siguiente:

```
sub {
    $_[0] =~ s/\A\s*//;
    $_[0] =~ s/\A(\d+(\.\d+)?)// and return ("$_", 'NUMBER');
    $_[0] =~ s/\A(#.*)// and return ("$_", 'COMMENT');
    $_[0] =~ s/\A("[^"]*" )// and return ("$_", 'STRING');
    $_[0] =~ s/\A(\$[a-zA-Z_]\w*)// and return ("$_", 'VAR');
    $_[0] =~ s/\A([a-zA-Z_]\w*)// and return ("$_", 'ID');
    $_[0] =~ s/\A(.)// and return ("$_", 'OTHER');
    $_[0] =~ s/(.|\n)// and return ("$_", ''); # 'tragamos' el resto
    return; # la cadena es vacía
}
```

Recordemos que el ancla `\A` casa con el comienzo de la cadena, incluso si esta es multilínea. Esta metodología es descrita por Conway en [?]. Obsérvese que una vez mas, Perl se aparta de la ortodoxia:

en otros lenguajes un objeto suele adoptar la forma de un `struct` o `hash`. En Perl cualquier cosa que se bendiga puede ser promocionada a la categoría de objeto. En este caso es una subrutina.

Veamos primero el código del constructor, situado en el módulo `Lexer.pm`:

```
package Lexer;
$VERSION = 1.00;
use strict;

sub new {
    my $class = shift; # El nombre del paquete
    my $spaces = shift; # Los espacios en blanco a "tragar"
    my $end_sub =<<EOS;
    \$_[0] =~ s/(.|\n)// and return ("\$1", ''); # 'tragamos' el resto
    return; # la cadena es vacía
}
EOS
    my $code = "sub {\n    \$_[0] =~ s/\\A$spaces//;\n";

    while (my ($regexp, $token) = splice @_, 0, 2) {
        my $lexer_line =<<EOL;
        \$_[0] =~ s/\\A($regexp)// and return ("\$1", '$token');
    EOL

        $code .= $lexer_line;
    }
    $code .= $end_sub;

    my $sub = eval $code or die "Error en la definición de los terminales\n";
    bless $sub, $class;
}
```

Obsérvese que, al bendecirla, estamos elevando la subrutina `$sub` a la categoría de objeto, facilitando de este modo

- la coexistencia de diversos analizadores léxicos en el mismo programa,
- el adjudicarle métodos al objeto y
- el uso de la extensión de la sintaxis flecha para objetos.

Así podemos tener diversos métodos asociados con el analizador léxico que modifiquen su conducta por defecto. Por ejemplo, podemos disponer de un método `extract_next` que recibiendo como entrada la cadena a analizar elimine de la misma el siguiente terminal:

```
($val, $token) = $lex->extract_next($input);
```

El código de dicho método será:

```
sub extract_next {
    &{$_[0]}($_[1]); # $lex->extract_next($data) <=> $lex($data)
}
```

pero es posible tener también un método `lookahead` que devuelva la misma información sin eliminar el terminal (aunque en el siguiente código si que se eliminan los blancos iniciales) de la cadena de entrada:

```

sub lookahead {
    my ($val, $token) = &{$_[0]}($_[1]);
    $_[1] = $val.$_[1];
    return ($val, $token);
}

```

y así, cuantos métodos hagan falta. Por ejemplo, podemos introducir un método `extract_all` que produzca la lista completa de parejas (cadena, identificador-de-terminal) sin destruir la entrada:

```

sub extract_all {
    my ($self, $input) = @_; # no destructivo
    my @tokens = ();

    while ($input) {
        push @tokens, &{$self}($input);
    }

    return @tokens;
}

```

Veamos un ejemplo de uso. El programa `uselexer.pl` hace uso del package/clase, creando primero un objeto de la clase `Lexer` y utilizando después dos de sus métodos:

```

#!/usr/bin/perl -w
use Lexer;

my ($val, $token);
my $lex = Lexer->new('\s*', #espacios
                    '\d+(\.\d+)?'=>'NUMBER',
                    '#.*'=>'COMMENT',
                    '"[^"]*"'=> 'STRING',
                    '\$[a-zA-Z_]\w*'=>'VAR',
                    '[a-zA-Z_]\w*'=>'ID',
                    '.*'=>'OTHER');

undef($/);

#token a token; lectura destructiva
my $input = <>;
while (($val, $token) = $lex->extract_next($input)) {
    print "$token -> $val\n";
}

# todos a la vez
$input = <>;
my @tokens = $lex->extract_all($input);
print "@tokens\n";

```

Ejemplo de ejecución:

```

$ ./uselexer.pl
$a = 34; # comentario

```

pulsamos CTRL-D y obtenemos la salida:

```

VAR -> $a
OTHER -> =
NUMBER -> 34

```

```
OTHER -> ;
COMMENT -> # comentario
```

De nuevo damos la misma entrada, pero ahora utilizaremos el método no destructivo `extract_all`:

```
$a = 34; #comentario
```

pulsamos CTRL-D de nuevo y obtenemos la salida:

```
$a VAR = OTHER 34 NUMBER ; OTHER #comentario COMMENT
```

5.2. Construcción usando la opción `g` y el ancla `G`

Tomemos como ejemplo el análisis léxico de una sencilla calculadora que admite constantes numéricas enteras, los paréntesis `(,)` y los operadores `+, -, * y /`. Se trata de producir a partir de la cadena de entrada una lista con los *tokens* y sus valores. Por ejemplo, la cadena `3*4-2` generará la lista: `('num', 3, 'op', '*', 'num', 4, 'op', '-', 'num', 2)`.

Las versiones mas recientes de Perl disponen de una opción `/c` que afecta a las operaciones de emparejamiento con `/g` en el contexto escalar. Normalmente, cuando una búsqueda global escalar tiene lugar y no ocurre casamiento, la posición `\G` es reestablecida al comienzo de la cadena. La opción `/c` hace que la posición inicial de emparejamiento permanezca donde la dejó el último emparejamiento con éxito. combinando `\G` y `/c` es posible escribir con cierta comodidad un analizador léxico:

```
{ # Con el redo del final hacemos un bucle "infinito"
  if (\G(\d+)/gc) {
    push @tokens, 'num', $1;
  } elsif (m|\G([-+*()])|gc) {
    push @tokens, 'pun', $1;
  } elsif (/G./gc) {
    die 'Caracter invalido';
  }
  else {
    last;
  }
  redo;
}
```

5.3. La clase `Parse::Lex`

La clase `Parse::Lex` nos permite crear un analizador léxico. La estrategia seguida es mover el puntero de búsqueda dentro de la cadena a analizar utilizando conjuntamente el operador `pos()` y el ancla `\G`.

```
> cat -n tokenizer.pl
1  #!/usr/local/bin/perl
2
3  require 5.004;
4  #BEGIN { unshift @INC, "../lib"; }
5
6  $^W = 0;
7  use Parse::Lex;
8  print STDERR "Version $Parse::ALex::VERSION\n";
9
10 @token = (
11          qw(
```



```

12         ADDOP      [-+]
13         LEFTP      [\([
14         RIGHTP     [\)]
15         INTEGER     [1-9][0-9]*
16         NEWLINE     \n
17     ),
18     qw(String),      [qw(" (?:[^"]+|"")* " )],
19     qw(ERROR .*), sub {
20         die qq!can\'t analyze: "$_[1]"!;
21     }
22 );
23
24 Parse::Lex->trace;
25 $lexer = Parse::Lex->new(@token);
26
27 $lexer->from(\*DATA);
28 print "Tokenization of DATA:\n";
29
30 TOKEN:while (1) {
31     $token = $lexer->next;
32     if (not $lexer->eoi) {
33         print "Record number: ", $lexer->line, "\n";
34         print "Type: ", $token->name, "\t";
35         print "Content:->", $token->text, "<-\n";
36     } else {
37         last TOKEN;
38     }
39 }
40
41 __END__
42 1+2-5
43 "This is a multiline
44 string with an embedded "" in it"
45 this is an invalid string with a "" in it"

```

La línea 3 usa la palabra reservada **require** la cual normalmente indica la existencia de una dependencia. Cuando se usa como:

```
require "file.pl";
```

simplemente carga el fichero en nuestro programa. Si se omite el sufijo, se asume **.pm**. Si el argumento, como es el caso, es un número o un número de versión, se compara el número de versión de perl (como se guarda en \$]) y se compara con el argumento. si es mas pequeño, se aborta la ejecución.

La línea 4, aunque comentada, contiene una definición del inicializador **BEGIN**. Si aparece, este código será ejecutado tan pronto como sea posible. el intérprete Perl busca los módulos en uno de varios directorios estándar contenidos en la variable de entorno **PERL5LIB**. Esa lista esta disponible en un programa Perl a través de la variable **@INC**. Recuerda que el compilador sustituye cada **::** por el separador de caminos. Asi la orden: **unshift @INC, "../lib"**; coloca este directorio entre los que Perl realiza la búsqueda por módulos. En mi máquina los módulos de análisis léxico están en:

```

> locate Lex.pm
/usr/local/lib/perl5/site_perl/5.8.0/Parse/ALex.pm
/usr/local/lib/perl5/site_perl/5.8.0/Parse/CLex.pm
/usr/local/lib/perl5/site_perl/5.8.0/Parse/Lex.pm
/usr/local/lib/perl5/site_perl/5.8.0/Parse/YYLex.pm

```

La línea 6 establece a 0 la variable `$^W` o `$WARNING` la cual guarda el valor actual de la opción de *warning*, si es cierto o falso.

La línea 7 indica que hacemos uso de la clase `Parse::Lex` y la línea 8 muestra el número de versión. La clase `ALex` es la clase abstracta desde la que heredan las otras. Si editas el módulo `ALex.pm` encontrarás el siguiente comentario que aclara la jerarquía de clases:

```
# Architecture:
#           Parse::ALex - Abstract Lexer
#           |
#           +-----+
#           |       |
#           |       Parse::Tokenizer
#           |       |       |
#           LexEvent  Lex  CLex  ...      - Concrete lexers
```

en ese mismo fichero encontrarás la declaración de la variable:

```
$Parse::ALex::VERSION = '2.15';
```

Las líneas de la 10 a la 22 definen el *array* `@token`:

```
10 @token = (
11     qw(
12         ADDOP      [-+]
13         LEFTP      [\([
14         RIGHTP     [\)]
15         INTEGER    [1-9][0-9]*
16         NEWLINE    \n
17     ),
18     qw(String),    [qw(" (?:[^"]+|")* " )],
19     qw(ERROR .*), sub {
20         die qq!can't analyze: "$_[1]"!n!;
21     }
22 );
```

Esto da lugar a un *array* de la forma:

```
"ADDOP", "[-+]",
"LEFTP", "[([",
"RIGHTP", "[)]",
"INTEGER", "[1-9][0-9]*",
"NEWLINE", "\n",
"STRING", ARRAY(0x811f55c),
"ERROR", ".*", CODE(0x82534b8)
```

Observa que los elementos 11 y 14 son referencias. El primero a un *array* anónimo conteniendo la expresión regular `"(?:[^"]+)"`— para las `STRING` y el segundo es una referencia a la subrutina anónima que muestra el mensaje de error.

La línea 24 llama al método `trace` como método de la clase, el cual activa el modo “traza”. La activación del modo traza debe establecerse antes de la creación del analizador léxico. Es por esto que precede a la llamada `$lexer = Parse::Lex->new(@token)` en la línea 25. Se puede desactivar el modo traza sin mas que hacer una segunda llamada a este método. El método admite la forma `trace OUTPUT`, siendo `OUTPUT` un nombre de fichero o un manipulador de ficheros. En este caso la traza se redirige a dicho fichero.

La llamada al método `from` en la línea 27 establece la fuente de los datos a ser analizados. el argumento de este método puede ser una cadena o una lista de cadenas o una referencia a un manipulador

de fichero. En este caso la llamada `$lexer->from(*DATA)` establece que la entrada se produce desde el manipulador de ficheros especial `DATA`, el cual se refiere a todo el texto que sigue al *token* `__END__` en el fichero que contiene el gui3n Perl.

La llamada al m3todo `next` en la l3nea 31 fuerza la b3squeda por el siguiente *token*. Devuelve un objeto de la clase `Parse::Token`. Existe un *token* especial, `Token::EOI`— que indica el final de los datos. En el ejemplo, el final se detecta a trav3s del m3todo `eoi`.

El m3todo `line` es usado para devolver el n3mero de l3nea del registro actual.

Los objetos *token* son definidos mediante la clase `Parse::Token`, la cual viene con `Parse::Lex`. la definici3n de un *token* normalmente conlleva un nombre simb3lico (por ejemplo, `INTEGER`), seguido de una expresi3n regular. Se puede dar como tercer argumento una referencia a una subrutina an3nima, en cuyo caso la subrutina es llamada cada vez que se reconoce el *token*. Los argumentos que recibe la subrutina son la referencia al objeto *token* y la cadena reconocida por la expresi3n regular. El valor retornado por la subrutina es usado como los nuevos contenidos de la cadena asociada con el objeto *token* en cuesti3n. Un objeto *token* tiene diversos m3todos asociados. Asi el m3todo `name` devuelve el nombre del *token*. El m3todo `text` devuelve la cadena de caracteres reconocida por medio del *token*. Se le puede pasar un par3metro `text` `EXPR` en cuyo caso `EXPR` define la cadena de caracteres asociada con el lexema.

```
> tokenizer.pl
Version 2.15
Trace is ON in class Parse::Lex
Tokenization of DATA:
[main::lexer|Parse::Lex] Token read (INTEGER, [1-9][0-9]*): 1
Record number: 1
Type: INTEGER    Content:->1<-
[main::lexer|Parse::Lex] Token read (ADDOP, [-+]): +
Record number: 1
Type: ADDOP      Content:->+<-
[main::lexer|Parse::Lex] Token read (INTEGER, [1-9][0-9]*): 2
Record number: 1
Type: INTEGER    Content:->2<-
[main::lexer|Parse::Lex] Token read (ADDOP, [-+]): -
Record number: 1
Type: ADDOP      Content:->-<-
[main::lexer|Parse::Lex] Token read (INTEGER, [1-9][0-9]*): 5
Record number: 1
Type: INTEGER    Content:->5<-
[main::lexer|Parse::Lex] Token read (NEWLINE, \n):

Record number: 1
Type: NEWLINE    Content:->
<-
[main::lexer|Parse::Lex] Token read (STRING, \"(?:[^\"]+|\"\\\" )*\"): "This is a multiline
string with an embedded "" in it"
Record number: 3
Type: STRING      Content:->"This is a multiline
string with an embedded "" in it"<-
[main::lexer|Parse::Lex] Token read (NEWLINE, \n):

Record number: 3
Type: NEWLINE    Content:->
<-
[main::lexer|Parse::Lex] Token read (ERROR, .*): this is an invalid string with a "" in it"
can't analyze: "this is an invalid string with a "" in it"
```

5.3.1. Condiciones de arranque

Los *tokens* pueden ser prefijados por “condiciones de arranque” o “estados”. Por ejemplo:

```
qw(C1:T1 er1), sub { # acción },
qw(T2 er2), sub { # acción },
```

el símbolo T1 será reconocido únicamente si el estado T1 está activo. Las condiciones de arranque se activan utilizando el método `start(condicion)` y se desactivan mediante `end(condicion)`. La llamada `start('INITIAL')` nos devuelve a la condición inicial.

```
> cat -n expectfloats.pl
1  #!/usr/bin/perl -w
2  use Parse::Lex;
3  @token = (
4      'EXPECT', 'expect-floats', sub {
5          $lexer->start('expect');
6          $_[1]
7      },
8      'expect:FLOAT', '\d+\.\d+',
9      'expect:NEWLINE', '\n', sub { $lexer->end('expect') ; $_[1] },
10     'expect:SIGN', '[+-]',
11     'NEWLINE2', '\n',
12     'INT', '\d+',
13     'DOT', '\.',
14     'SIGN2', '[+-]'
15 );
16
17 Parse::Lex->exclusive('expect');
18 $lexer = Parse::Lex->new(@token);
19
20 $lexer->from(\*DATA);
21
22 TOKEN:while (1) {
23     $token = $lexer->next;
24     if (not $lexer->eoi) {
25         print $token->name," ";
26         print "\n" if ($token->text eq "\n");
27     } else {
28         last TOKEN;
29     }
30 }
31
32 __END__
33 1.4+2-5
34 expect-floats 1.4+2.3-5.9
35 1.5
```

Ejecución:

```
> expectfloats.pl
INT DOT INT SIGN2 INT SIGN2 INT NEWLINE2
EXPECT FLOAT SIGN FLOAT SIGN FLOAT NEWLINE
INT DOT INT NEWLINE2
```

El siguiente ejemplo elimina los comentarios anidados en un programa C:

```

> cat -n nestedcom.pl
1  #!/usr/local/bin/perl -w
2
3  require 5.004;
4  #BEGIN { unshift @INC, "../lib"; }
5
6  $^W = 0;
7  use Parse::Lex;
8
9  @token = (
10         'STRING',          '"([^\"]|\\.)"', sub { print "$_[1]"; $_[1]; },
11         'CI',               '\\/*', sub { $lexer->start('comment'); $c++; $_[1]; },
12         'CF',               '\\*\\/', sub { die "Error, comentario no finalizado!"; },
13         'OTHER',            '(.|\\n)', sub { print "$_[1]"; $_[1]; },
14         'comment:CCI',      '\\/*', sub { $c++; $_[1]; },
15         'comment:CCF',      '\\*\\/', sub { $c--; $lexer->end('comment') if ($c == 0); },
16         'comment:COTHER',   '(.|\\n)',
17     );
18
19     #Parse::Lex->trace;
20     Parse::Lex->exclusive('comment');
21     Parse::Lex->skip('');
22
23     $lexer = Parse::Lex->new(@token);
24
25     $lexer->from(\\*DATA);
26
27     $lexer = Parse::Lex->new(@token);
28
29     $lexer->from(\\*DATA);
30
31     TOKEN:while (1) {
32         $token = $lexer->next;
33         last TOKEN if ($lexer->eoi);
34     }
35
36     __END__
37     main() { /* comment */
38         printf("hi! /* \\\"this\\\" is not a comment */"); /* another /*nested*/
39             comment */
40     }

```

Al ejecutarlo, produce la salida:

```

> nestedcom.pl
main() {
    printf("hi! /* \\\"this\\\" is not a comment */");
}

```

5.4. La Clase Parse::CLex

Relacionada con Parse::Lex está la clase Parse::CLex, la cual avanza consumiendo la cadena analizada mediante el uso del operador de sustitución (s///). Los analizadores producidos mediante

esta segunda clase no permiten el uso de anclas en las expresiones regulares. Tampoco disponen de acceso a la subclase `Parse::Token`. He aquí el mismo ejemplo, usando la clase `Parse::CLex`:

```
> cat -n ctokenizer.pl
1  #!/usr/local/bin/perl -w
2
3  require 5.000;
4  BEGIN { unshift @INC, "../lib"; }
5  use Parse::CLex;
6
7  @token = (
8      qw(
9          ADDOP      [-+]
10         LEFTP      [\([
11         RIGHTP     [\)]
12         INTEGER    [1-9][0-9]*
13         NEWLINE    \n
14     ),
15     qw(String),    [qw(" (?:[^"]+|")* ")],
16     qw(ERROR      .*), sub {
17         die qq!can't analyze: "$_[1]"!;
18     }
19 );
20
21 Parse::CLex->trace;
22 $lexer = Parse::CLex->new(@token);
23
24 $lexer->from(\*DATA);
25 print "Tokenization of DATA:\n";
26
27 TOKEN:while (1) {
28     $token = $lexer->next;
29     if (not $lexer->eoi) {
30         print "Record number: ", $lexer->line, "\n";
31         print "Type: ", $token->name, "\t";
32         print "Content:->", $token->getText, "<-\n";
33     } else {
34         last TOKEN;
35     }
36 }
37
38 __END__
39 1+2-5
40 "This is a multiline
41 string with an embedded "" in it"
42 this is an invalid string with a "" in it"
43
44

> ctokenizer.pl
Trace is ON in class Parse::CLex
Tokenization of DATA:
[main::lexer|Parse::CLex] Token read (INTEGER, [1-9][0-9]*): 1
Record number: 1
```

```

Type: INTEGER    Content:->1<-
[main::lexer|Parse::CLex] Token read (ADDOP, [-+]): +
Record number: 1
Type: ADDOP      Content:->+<-
[main::lexer|Parse::CLex] Token read (INTEGER, [1-9][0-9]*): 2
Record number: 1
Type: INTEGER    Content:->2<-
[main::lexer|Parse::CLex] Token read (ADDOP, [-+]): -
Record number: 1
Type: ADDOP      Content:->-<-
[main::lexer|Parse::CLex] Token read (INTEGER, [1-9][0-9]*): 5
Record number: 1
Type: INTEGER    Content:->5<-
[main::lexer|Parse::CLex] Token read (NEWLINE, \n):

Record number: 1
Type: NEWLINE    Content:->
<-
[main::lexer|Parse::CLex] Token read (STRING, \"(?:[^\"]+|\\\"\\\")*\"): "This is a multiline
string with an embedded "" in it"
Record number: 3
Type: STRING      Content:->"This is a multiline
string with an embedded "" in it"<-
[main::lexer|Parse::CLex] Token read (NEWLINE, \n):

Record number: 3
Type: NEWLINE    Content:->
<-
[main::lexer|Parse::CLex] Token read (ERROR, .*): this is an invalid string with a "" in it"
can't analyze: "this is an invalid string with a "" in it" at ctokenizer.pl line 17, <DATA> 1

```

5.5. El Módulo Parse::Flex

```

lhp@nereida:/tmp/Parse-Flex-0.11/t$ cat -n cas.1
 1  %{
 2  #define YY_DECL char* yylex(void)
 3  %}
 4
 5  %pointer
 6  %option yylineno noyywrap
 7
 8  NUM    [0-9]+
 9  WORD   [a-zA-Z][A-Za-z0-9_]*
10  EMAIL  ({WORD}[.]?)+[@]({WORD}.)+[a-z]{2,3}
11
12  %%
13  {EMAIL}      return "EMAIL" ;
14  {NUM}        return "NUM"   ;
15  {WORD}       return "WORD"  ;
16  [a-zA-Z]/\   return "CHAR"  ;
17  \n
18  [\t ]
19  .

```

```

20 <<EOF>>      return ""      ;
21 %%
22
lhp@nereida:/tmp/Parse-Flex-0.11/t$ makelexer.pl cas.l
lhp@nereida:/tmp/Parse-Flex-0.11/t$ ls -ltr | tail -2
-rw-r--r--  1 lhp lhp  1933 2006-06-28 19:07 Flexer28424.pm
-rwxr-xr-x  1 lhp lhp 38077 2006-06-28 19:07 Flexer28424.so
lhp@nereida:/tmp/Parse-Flex-0.11/t$ perl -MFlexer28424 -de 0
Loading DB routines from perl5db.pl version 1.28
main::(-e:1):  0
    DB<1> $x = gen_walker()
    DB<2> x $x->()
tut@titi.top
0  'EMAIL'
1  'tut@titi.top'
    DB<3> x $x->()
funchal
0  'WORD'
1  'funchal'
    DB<4> x $x->()
432
0  'NUM'
1  432
    DB<5> x $x->()
;
    .:::
fin
0  'WORD'
1  'fin'
    DB<6> q

```

5.6. Usando Text::Balanced

La función `extract_multiple` del módulo `Text::Balanced` puede ser considerada un generador de analizadores léxicos. El primer argumento es la cadena a ser procesada y el segundo la lista de *extractores* a aplicar a dicha cadena. Un extractor puede ser una subrutina, pero también una expresión regular o una cadena.

En un contexto de lista devuelve la lista con las subcadenas de la cadena original, según han sido producidas por los extractores. En un contexto escalar devuelve la primera subcadena que pudo ser extraída de la cadena original. En cualquier caso, `extract_multiple` empieza en la posición actual de búsqueda indicada por `pos` y actualiza el valor de `pos` de manera adecuada.

Si el extractor es una referencia a un hash, deberá contener exactamente un elemento. La clave actúa como nombre del terminal o *token* y el valor es un extractor. La clave es utilizada para bendecir la cadena casada por el extractor.

Veamos un programa de ejemplo:

```

$ cat extract_variable.pl
#!/usr/local/bin/perl5.8.0 -w

use strict;
use Text::Balanced qw( extract_variable extract_quotelike
                      extract_codeblock extract_multiple);

my $text = <<'EOS'

```



```
#!/usr/local/bin/perl5.8.0 -w

$pattern = shift || '';
$pattern =~ s{::}{/}g;
$pattern =~ s{$}{.pm};
@dirs = qw(
/usr/local/lib/perl5/5.8.0/i686-linux
/usr/local/lib/perl5/5.8.0
/usr/local/lib/perl5/site_perl/5.8.0/i686-linux
/usr/local/lib/perl5/site_perl/5.8.0
/usr/local/lib/perl5/site_perl
/usr/lib/perl5
/usr/share/perl5
/usr/lib/perl/5.6.1
/usr/share/perl/5.6.1
);

for (@dirs) {
    my $file = $_."/".$pattern;
    print "$file\n" if (-e $file);
}

EOS
;
my @tokens = &extract_multiple(
    $text,
    [
        {variable => sub { extract_variable($_[0], '') }},
        {quote => sub { extract_quotelike($_[0], '') }},
        {code => sub { extract_codeblock($_[0], '{}', '') }},
        {comment => qr/#.*/}
    ],
    undef,
    0
)
;

for (@tokens) {
    if (ref($_)) {
        print ref($_), "=> $_\n";
    }
    else {
        print "other things: $_\n";
    }
}
}
```

La subrutina `extract_variable` extrae cualquier variable Perl válida. El primer argumento es la cadena a ser procesada, el segundo es una cadena que especifica un patrón indicando el prefijo a saltarse. Si se omite como en el ejemplo, se usarán blancos. La subrutina `extract_quotelike` extrae cadenas Perl en cualquiera de sus múltiples formas de representación. La subrutina `extract_codeblock` extrae un bloque de código.

Ejecución del programa anterior:

```
$ ./extract_variable.pl
```

```
comment=> #!/usr/local/bin/perl5.8.0 -w
other things:
```

```
variable=> $pattern
other things: = shift ||
quote=> ''
other things: ;
```

```
variable=> $pattern
other things: =~
quote=> s{::}{/}g
other things: ;
```

```
variable=> $pattern
other things: =~
quote=> s{$}{.pm}
other things: ;
```

```
variable=> @dirs
other things: =
quote=> qw(
/usr/local/lib/perl5/5.8.0/i686-linux
/usr/local/lib/perl5/5.8.0
/usr/local/lib/perl5/site_perl/5.8.0/i686-linux
/usr/local/lib/perl5/site_perl/5.8.0
/usr/local/lib/perl5/site_perl
/usr/lib/perl5
/usr/share/perl5
/usr/lib/perl/5.6.1
/usr/share/perl/5.6.1
)
other things: ;
```

```
for (
variable=> @dirs
other things: )
code=> {
    my $file = $_."/".$pattern;
    print "$file\n" if (-e $file);
}
other things:
```

Capítulo 6

RecDescent

6.1. Introducción

`Parse::RecDescent` es un módulo, escrito por Damian Conway, que permite construir analizadores sintácticos descendentes con vuelta atrás. No forma parte de la distribución estándar de Perl, por lo que es necesario descargarlo desde CPAN (consulte la sección [?] para una descripción de como hacerlo). En lo que sigue, y por brevedad, nos referiremos a el como RD.

```
1 #!/usr/local/bin/perl5.8.0 -w
2 use strict;
3 use warnings;
4
5 use Parse::RecDescent;
6
7
8 my $grammar = q {
9
10  start:  seq_1 seq_2
11
12  seq_1   :   'A' 'B' 'C' 'D'
13           { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
14           | 'A' 'B' 'C' 'D' 'E' 'F'
15           { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
16
17  seq_2   : character(s)
18
19  character: /\w/
20           { print "character: $item[1]\n" }
21
22 };
23
24 my $parser=Parse::RecDescent->new($grammar);
25
26 $parser->start("A B C D E F");
```

En la línea 5 se declara el uso del módulo. En la línea 8 se guarda la gramática en la variable escalar `$grammar`. La variable sintáctica `start` es considerada de arranque de la gramática. Las partes derechas de las producciones se separan mediante la barra |, y las acciones se insertan entre llaves. El contenido de las acciones es código Perl. Dentro de esas acciones es posible hacer alusión a los atributos de los símbolos de la producción a través del array `@item`. Así `$item[1]` contendrá el atributo asociado al primer símbolo (el carácter A en este caso). Para que una producción tenga éxito es necesario que la acción devuelva un valor definido (no necesariamente verdadero). El valor asociado con una producción

es el valor retornado por la acción, o bien el valor asociado con la variable `$return` (si aparece en la acción) o con el ítem asociado con el último símbolo que casó con éxito.

Es posible indicar una o mas repeticiones de una forma sentencial como en la línea 17. La regla:

```
seq_2 : character(s)
```

indica que el lenguaje `seq_2` se forma por la repetición de uno o mas caracteres (`character`). El valor asociado con una pluralidad de este tipo es una referencia a un array conteniendo los valores individuales de cada uno de los elementos que la componen, en este caso de los `character`. Es posible introducir un *patrón de separación* que indica la forma en la que se deben separar los elementos:

```
seq_2 : character(s /;/)
```

lo cual especifica que entre `character` y `character` el analizador debe saltarse un punto y coma.

Como se ve en la línea 19 es legal usar expresiones regulares en la parte derecha. La llamada al constructor en la línea 24 deja en `$parser` un analizador sintáctico recursivo descendente. El analizador así construido dispone de un método por cada una de las variables sintácticas de la gramática. El método asociado con una variable sintáctica reconoce el lenguaje generado por esa variable sintáctica. Así, la llamada de la línea 26 al método `start` reconoce el lenguaje generado por `start`. Si la variable sintáctica tiene éxito el método retorna el valor asociado. En caso contrario se devuelve `undef`. Esto hace posible retornar valores como 0, "0", or "".

La siguiente ejecución muestra el orden de recorrido de las producciones.

```
$ ./firstprodnolongest.pl
seq_1: A B C D
character: E
character: F
```

¿Que ha ocurrido? RD es perezoso en la evaluación de las producciones. Como la primera regla tiene éxito y es evaluada antes que la segunda, retorna el control a la parte derecha de la regla de `start`. A continuación se llama a `seq2`.

En este otro ejemplo invertimos el orden de las reglas de producción de `seq_1`:

```
#!/usr/local/bin/perl5.8.0 -w

use strict;
use warnings;

use Parse::RecDescent;

my $grammar = q {

    start:  seq_1 seq_2

    seq_1   :  'A' 'B' 'C' 'D' 'E' 'F'
              { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
    |  'A' 'B' 'C' 'D'
              { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }

    seq_2   :  character(s)

    character: /\w/
              { print "character: $item[1]\n" }
```

```
};

my $parser=Parse::RecDescent->new($grammar);
```

```
$parser->start("A B C D E F");
```

La ejecución resultante da lugar al éxito de la primera regla, que absorbe toda la entrada:

```
$ ./firstprodnolongest2.pl
seq_1: A B C D E F
```

6.2. Orden de Recorrido del Árbol de Análisis Sintáctico

RD es un analizador sintáctico descendente recursivo pero no es predictivo. Cuando estamos en la subrutina asociada con una variable sintáctica, va probando las reglas una a una hasta encontrar una que tiene éxito. En ese momento la rutina termina retornando el valor asociado. Consideremos la siguiente gramática (véase el clásico libro de Aho et al. [?], página 182):

$$S \rightarrow cAd$$

$$A \rightarrow ab \mid a$$

Claramente la gramática no es LL(1). Codifiquemos la gramática usando RD:

```
$ cat -n aho182.pl
1  #!/usr/bin/perl -w
2  use strict;
3  use warnings;
4
5  use Parse::RecDescent;
6
7  $::RD_TRACE = 1;
8  $::RD_AUTOACTION = q{ 1 };
9  my $grammar = q {
10   S : 'c' A 'd' { print "Éxito!\n" }
11   A : 'a' 'b'
12       | 'a'
13 };
14
15 my $parser=Parse::RecDescent->new($grammar);
16 my $input = <>;
17 $parser->S($input);
```

en este ejemplo usamos dos nuevas variables (líneas 7 y 8) que gobiernan la conducta de un analizador generado por RD. Las variables `$::RD_TRACE` y `$::RD_AUTOACTION`. La primera hace que el analizador generado vuelque en `stderr` un informe del análisis. La asignación en la línea 8 de una cadena a `$::RD_AUTOACTION` hace que todas las producciones con las cuáles no se asocie una acción explícitamente, tengan como acción asociada la cadena que figura en la variable `$::RD_AUTOACTION` (en este caso, devolver 1). Así las producciones en las líneas 11 y 12 tienen asociadas la acción `{ 1 }`, mientras que la acción asociada con la producción en la línea 10 es `{ print "Éxito!\n" }`.

Obsérvese que, dado que el analizador se ejecuta dentro de su propio espacio de nombres, las variables que pertenezcan al paquete `Main` como es el caso de `RD_AUTOACTION` deben ser explicitadas prefijándolas con el nombre del paquete.

Veamos el comportamiento del programa con la entrada `c a d`. Al ejecutar el programa, lo primero que ocurre, al construir el objeto analizador en la línea 15 es que RD nos informa de su interpretación de los diferentes símbolos en la gramática:

```

$ ./aho182.pl
Parse::RecDescent: Treating "S :" as a rule declaration
Parse::RecDescent: Treating "c" as a literal terminal
Parse::RecDescent: Treating "A" as a subrule match
Parse::RecDescent: Treating "d" as a literal terminal
Parse::RecDescent: Treating "{ print "Éxito!\n" }" as an action
Parse::RecDescent: Treating "A :" as a rule declaration
Parse::RecDescent: Treating "a" as a literal terminal
Parse::RecDescent: Treating "b" as a literal terminal
Parse::RecDescent: Treating "|" as a new production
Parse::RecDescent: Treating "a" as a literal terminal
printing code (12078) to RD_TRACE

```

Ahora se ejecuta la entrada en la línea 16, tecleamos `c a d` y sigue una información detallada de la construcción del árbol sintáctico concreto:

```

c a d
1|  S      |Trying rule: [S]                |
1|  S      |                                |"c a d\n"
1|  S      |Trying production: ['c' A 'd']           |
1|  S      |Trying terminal: ['c']                    |
1|  S      |>>Matched terminal<< (return value:      |
|          |[c])                          |

```

La salida aparece en cuatro columnas. El número de línea en el texto de la gramática, la variable sintáctica (que en la jerga RD se denomina regla o subregla), una descripción verbal de lo que se está haciendo y por último, la cadena de entrada que queda por leer. Se acaba de casar `c` y por tanto se pasa a llamar al método asociado con `A`. en la entrada queda `a d\n`:

```

1|  S      |                                |" a d\n"
1|  S      |Trying subrule: [A]                |
2|  A      |Trying rule: [A]                   |
2|  A      |Trying production: ['a' 'b']        |
2|  A      |Trying terminal: ['a']              |
2|  A      |>>Matched terminal<< (return value:  |
|          |[a])                    |
2|  A      |                                |" d\n"
2|  A      |Trying terminal: ['b']              |
2|  A      |<<Didn't match terminal>>          |

```

RD intenta infructuosamente la primera producción

$$A \rightarrow a b$$

por tanto va a probar con la segunda producción de `A`:

```

2|  A      |                                |"d\n"
2|  A      |Trying production: ['a']            |
2|  A      |                                |" a d\n"
2|  A      |Trying terminal: ['a']              |
2|  A      |>>Matched terminal<< (return value:  |
|          |[a])                    |
2|  A      |                                |" d\n"
2|  A      |Trying action                      |
2|  A      |>>Matched action<< (return value: [1])|
2|  A      |>>Matched production: ['a']<<      |

```

```

2|   A   |>>Matched rule<< (return value: [1]) |
2|   A   |(consumed: [ a]) |
1|   S   |>>Matched subrule: [A]<< (return |
|       |value: [1] |
1|   S   |Trying terminal: ['d'] |
1|   S   |>>Matched terminal<< (return value: |
|       |[d]) |
1|   S   | | "\n"
1|   S   |Trying action |
Éxito!
1|   S   |>>Matched action<< (return value: [1])|
1|   S   |>>Matched production: ['c' A 'd']<< |
1|   S   |>>Matched rule<< (return value: [1]) |
1|   S   |(consumed: [c a d]) |

```

De este modo la frase es aceptada. Sin embargo, ¿que ocurriría si invertimos el orden de las producciones de A y damos como entrada la cadena *c a b d*? Puesto que RD es ocioso, la ahora primera regla A : a tendrá éxito y el control retornará al método asociado con S:

```

$ ./aho182_reverse.pl
c a b d
1|   S   |Trying rule: [S] |
1|   S   | | "c a b d\n"
1|   S   |Trying production: ['c' A 'd'] |
1|   S   |Trying terminal: ['c'] |
1|   S   |>>Matched terminal<< (return value: |
|       |[c]) |
1|   S   | | " a b d\n"
1|   S   |Trying subrule: [A] |
2|   A   |Trying rule: [A] |
2|   A   |Trying production: ['a'] |
2|   A   |Trying terminal: ['a'] |
2|   A   |>>Matched terminal<< (return value: |
|       |[a]) |
2|   A   | | " b d\n"
2|   A   |Trying action |
2|   A   |>>Matched action<< (return value: [1])|
2|   A   |>>Matched production: ['a']<< |
2|   A   |>>Matched rule<< (return value: [1]) |
2|   A   |(consumed: [ a]) |

```

Ahora RD va a rechazar la cadena. No intenta llamar de nuevo al método asociado con A para que haga nuevos intentos:

```

1|   S   |>>Matched subrule: [A]<< (return |
|       |value: [1] |
1|   S   |Trying terminal: ['d'] |
1|   S   |<<Didn't match terminal>> |
1|   S   | | "b d\n"
1|   S   |<<Didn't match rule>> |

```

Debes, por tanto, ser cuidadoso con el orden en el que escribes las producciones. El orden en que las escribas tiene una importante incidencia en la conducta del analizador. *En general el analizador generado por RD no acepta el lenguaje generado por la gramática*, en el sentido formal de la expresión *lenguaje generado*, tal y como fué definida en la definición 8.1.4. En este sentido, hay una separación

semántica similar a la que ocurre en las expresiones regulares. El *or* en las expresiones regulares de Perl también es perezoso, a diferencia de la convención adoptada en el estudio teórico de los lenguajes regulares y los autómatas.

Ejercicio 6.2.1. *Dado el siguiente programa Perl*

```
#!/usr/local/bin/perl5.8.0 -w
use strict;
use warnings;
use Parse::RecDescent;

$::RD_TRACE = 1;
my $grammar = q {

    start:  seq_1 ';'
    seq_1   : 'A' 'B' 'C' 'D'
              { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
            | 'A' 'B' 'C' 'D' 'E' 'F'
              { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
};

my $parser=Parse::RecDescent->new($grammar);
my $result = $parser->start("A B C D E F ;");
if (defined($result)) {
    print "Valida\n";
} else { print "No reconocida\n"; }
```

¿Cuál es la salida?

Indica que ocurre si cambiamos el lugar en el que aparece el separador punto y coma, poniéndolo en las producciones de seq_1 como sigue:

```
#!/usr/local/bin/perl5.8.0 -w
use strict;
use warnings;
use Parse::RecDescent;

$::RD_TRACE = 1;
my $grammar = q {

    start:  seq_1
    seq_1   : 'A' 'B' 'C' 'D' ';'
              { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
            | 'A' 'B' 'C' 'D' 'E' 'F' ';'
              { print "seq_1: " . join (" ", @item[1..$#item]) . "\n" }
};

my $parser=Parse::RecDescent->new($grammar);

my $result = $parser->start("A B C D E F ;");
if (defined($result)) {
    print "Valida\n";
} else { print "No reconocida\n"; }
```

¿Podrías resumir en palabras la forma en la que RD explora el árbol de análisis concreto? Intenta escribir un pseudocódigo que resuma el comportamiento de RD.

6.3. La ambigüedad de las sentencias if-then-else

Existen lenguajes que son *intrínsecamente ambiguos*, esto lenguajes independientes del contexto tales que, cualquier gramática que los genere es ambigua. Sin embargo, lo normal es que siempre se pueda encontrar una gramática que no sea ambigua y que genere el mismo lenguaje.

Existe una ambigüedad en la parte de las sentencias condicionales de un lenguaje. Por ejemplo, consideremos la gramática:

```
SENTENCIA → if EXPRESION then SENTENCIA
SENTENCIA → if EXPRESION then SENTENCIA else SENTENCIA
SENTENCIA → OTRASSENTENCIAS
```

Aquí OTRASSENTENCIAS es un terminal/comodín para aludir a cualesquiera otras sentencias que el lenguaje pueda tener (bucles, asignaciones, llamadas a subrutinas, etc.)

La gramática es ambigua, ya que para una sentencia como

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$$

existen dos árboles posibles: uno que asocia el “else” con el primer “if” y otra que lo asocia con el segundo. Los dos árboles corresponden a las dos posibles parentizaciones:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1 \text{ else } S_2)$$

Esta es la regla de prioridad usada en la mayor parte de los lenguajes: un “else” casa con el “if” mas cercano. la otra posible parentización es:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1) \text{ else } S_2$$

El siguiente ejemplo considera una aproximación a la resolución del problema usando `Parse::RecDescent`. Para simplificar, la única variable sintáctica que permanece del ejemplo anterior es SENTENCIA (que ha sido renombrada `st`). Las demás han sido convertidas en terminales y se han abreviado. En esta gramática colocamos primero la regla mas larga `st : 'iEt' st 'e' st`.

```
1 #!/usr/local/bin/perl5.8.0 -w
2 use strict;
3 use Parse::RecDescent;
4 use Data::Dumper;
5
6 $::RD_TRACE = 1;
7 $::RD_AUTOACTION = q{ [@item] };
8 my $grammar = q{
9   prog : st ';'
10   st : 'iEt' st 'e' st
11       | 'iEt' st
12       | 'o' { 'o' }
13 };
14
15 my $parse = Parse::RecDescent->new($grammar);
16
17 my $line;
18 while ($line = <>) {
19   print "$line\n";
20   my $result = $parse->prog($line);
21   if (defined($result)) { print Dumper($result); }
22   else { print "Cadena no válida\n"; }
23 }
```

En este caso, la variable `$::RD_AUTOACTION` ha sido establecida a `q{ [@item] }`. Esto significa que cada vez que una producción tenga éxito se devolverá una referencia al array conteniendo los atributos de los símbolos en la producción. El contenido de esta variable es evaluado siempre que la producción en cuestión no tenga una acción asociada explícitamente: Por tanto, se aplica para todas las reglas salvo para la regla en la línea 12, en la que se devuelve el carácter 'o'.

El programa anterior da lugar a la siguiente ejecución:

```
bash-2.05b$ ./ifthenelse.pl file4.txt
```

El fichero `file4.txt` contiene una sólo línea:

```
$ cat file4.txt
iEt iEt o e o ;
```

que es una frase con dos posibles árboles. RD intentará primero reiteradamente la primera de las producciones de `st`:

```
1| prog |Trying rule: [prog] |
1| prog | | |"iEt iEt o e o ;\n"
1| prog |Trying production: [st ';' ] |
1| prog |Trying subrule: [st] |
2| st |Trying rule: [st] |
2| st |Trying production: ['iEt' st 'e' st] |
2| st |Trying terminal: ['iEt'] |
2| st |>>Matched terminal<< (return value: |
| | | [iEt]) |
2| st | | |" iEt o e o ;\n"
2| st |Trying subrule: [st] |
3| st |Trying rule: [st] |
3| st |Trying production: ['iEt' st 'e' st] |
3| st |Trying terminal: ['iEt'] |
3| st |>>Matched terminal<< (return value: |
| | | [iEt]) |
3| st | | |" o e o ;\n"
```

Esta opción acabará fracasando, ya que sólo hay un `else`.

```
3| st |Trying subrule: [st] |
4| st |Trying rule: [st] |
4| st |Trying production: ['iEt' st 'e' st] |
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
```

Estamos delante del terminal `o` y obviamente la primera producción de `st` no casa, se intenta con la segunda, la cual naturalmente fracasará:

```
4| st | | |"o e o ;\n"
4| st |Trying production: ['iEt' st] |
4| st | | |" o e o ;\n"
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | |"o e o ;\n"
4| st |Trying production: ['o'] |
```

Seguimos delante del terminal `o` y la segunda producción de `st` tampoco casa, se intenta con la tercera producción:

```

4| st | | " o e o ;\n"
4| st |Trying terminal: ['o'] |
4| st |>>Matched terminal<< (return value: |
| | [o]) |
4| st | | " e o ;\n"
4| st |Trying action |
4| st |>>Matched action<< (return value: [o])|

```

La o por fin ha sido emparejada. Se ha ejecutado la acción no automática.

```

4| st |>>Matched production: ['o']<< |
4| st |>>Matched rule<< (return value: [o]) |
4| st |(consumed: [ o]) |
3| st |>>Matched subrule: [st]<< (return |
| | value: [o] |

```

Recuérdese por donde íbamos conjeturando. Se ha construido el árbol (erróneo):

```

st
|
' 'iEt' st 'e' st
|
' 'iEt' st 'e' st
| ^
' ,o'

```

Por ello, el e va a ser consumido sin problema:

```

3| st |Trying terminal: ['e'] |
3| st |>>Matched terminal<< (return value: |
| | [e]) |
3| st | | " o ;\n"
3| st |Trying subrule: [st] |
4| st |Trying rule: [st] |
4| st |Trying production: ['iEt' st 'e' st] |
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | "o ;\n"
4| st |Trying production: ['iEt' st] |
4| st | | " o ;\n"
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | "o ;\n"
4| st |Trying production: ['o'] |
4| st | | " o ;\n"
4| st |Trying terminal: ['o'] |
4| st |>>Matched terminal<< (return value: |
| | [o]) |
4| st | | " ;\n"
4| st |Trying action |
4| st |>>Matched action<< (return value: [o])|
4| st |>>Matched production: ['o']<< |
4| st |>>Matched rule<< (return value: [o]) |
4| st |(consumed: [ o]) |
3| st |>>Matched subrule: [st]<< (return |

```

```

|           |value: [o]
3|    st    |Trying action
3|    st    |>>Matched action<< (return value:
|           |[ARRAY(0x830cc04)])

```

Hemos construido el árbol:

```

st
|
' 'iEt' st 'e' st
      | ^
      ' 'iEt' st 'e' st
            | |
            ' 'o' 'o'

```

Así pues la segunda regla conjeturada `st : 'iEt' st 'e' st` ha tenido éxito. Después de esto esperamos ver `e`, pero lo que hay en la entrada es un punto y coma.

```

3|    st    |>>Matched production: ['iEt' st 'e'
|           |st]<<
3|    st    |>>Matched rule<< (return value:
|           |[ARRAY(0x830cc04)])
3|    st    |(consumed: [ iEt o e o])
2|    st    |>>Matched subrule: [st]<< (return
|           |value: [ARRAY(0x830cc04)])
2|    st    |Trying terminal: ['e']
2|    st    |<<Didn't match terminal>>
2|    st    |
|           |";\n"

```

Se ha fracasado. Se probará a este nivel con la siguiente regla `st : 'iEt' st` (que acabará produciendo un árbol de acuerdo con la regla del “if” mas cercano). Mágicamente, la entrada consumida es devuelta para ser procesada de nuevo (obsérvese la tercera columna). Sin embargo, los efectos laterales que las acciones ejecutadas por las acciones asociadas con los falsos éxitos permanecen. En este caso, por la forma en la que hemos escrito las acciones, no hay ningún efecto lateral.

```

2|    st    |Trying production: ['iEt' st]
2|    st    |
2|    st    |Trying terminal: ['iEt']
2|    st    |>>Matched terminal<< (return value:
|           |[iEt])
2|    st    |
2|    st    |" iEt o e o ;\n"
2|    st    |Trying subrule: [st]
3|    st    |Trying rule: [st]
3|    st    |Trying production: ['iEt' st 'e' st]

```

Ahora el analizador va a intentar el árbol:

```

st
|
' 'iEt' st
      |
      ' 'iEt' st 'e' st

```

que se corresponde con la interpretación clásica: *El else casa con el if mas cercano*.

```

3|    st    |Trying terminal: ['iEt']
3|    st    |>>Matched terminal<< (return value:
|           |[iEt])
3|    st    |
3|    st    |" o e o ;\n"
3|    st    |Trying subrule: [st]

```

La ineficiencia de RD es clara aqui. Va a intentar de nuevo las diferentes reglas hasta dar con la del 'o'

```

4| st |Trying rule: [st] |
4| st |Trying production: ['iEt' st 'e' st] |
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | "o e o ;\n"
4| st |Trying production: ['iEt' st] |
4| st | | " o e o ;\n"
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | "o e o ;\n"
4| st |Trying production: ['o'] |
4| st | | " o e o ;\n"
4| st |Trying terminal: ['o'] |
4| st |>>Matched terminal<< (return value: |
| | [o]) |
4| st | | " e o ;\n"
4| st |Trying action |
4| st |>>Matched action<< (return value: [o])|
4| st |>>Matched production: ['o']<< |
4| st |>>Matched rule<< (return value: [o]) |
4| st |(consumed: [ o]) |
3| st |>>Matched subrule: [st]<< (return |
| |value: [o] |

```

Ahora el "else" será consumido por la sentencia condicional anidada:

```

3| st |Trying terminal: ['e'] |
3| st |>>Matched terminal<< (return value: |
| | [e]) |
3| st | | " o ;\n"

```

y de nuevo debemos pasar por el calvario de todas las reglas, ya que la 'o' es la última de las reglas:

```

3| st |Trying subrule: [st] |
4| st |Trying rule: [st] |
4| st |Trying production: ['iEt' st 'e' st] |
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | "o ;\n"
4| st |Trying production: ['iEt' st] |
4| st | | " o ;\n"
4| st |Trying terminal: ['iEt'] |
4| st |<<Didn't match terminal>> |
4| st | | "o ;\n"
4| st |Trying production: ['o'] |
4| st | | " o ;\n"
4| st |Trying terminal: ['o'] |
4| st |>>Matched terminal<< (return value: |
| | [o]) |
4| st | | " ;\n"
4| st |Trying action |
4| st |>>Matched action<< (return value: [o])|

```

```

4|  st  |>>Matched production: ['o']<<      |
4|  st  |>>Matched rule<< (return value: [o]) |
4|  st  |(consumed: [ o])                    |

```

A partir de aquí todo encaja. Nótese la ejecución de la acción automática:

```

3|  st  |>>Matched subrule: [st]<< (return    |
|      |value: [o]                          |
3|  st  |Trying action                          |
3|  st  |>>Matched action<< (return value:     |
|      |[ARRAY(0x82f2774)])               |
3|  st  |>>Matched production: ['iEt' st 'e'   |
|      |st]<<                               |
3|  st  |>>Matched rule<< (return value:        |
|      |[ARRAY(0x82f2774)])               |
3|  st  |(consumed: [ iEt o e o])            |
2|  st  |>>Matched subrule: [st]<< (return      |
|      |value: [ARRAY(0x82f2774)])         |
2|  st  |Trying action                          |
2|  st  |>>Matched action<< (return value:     |
|      |[ARRAY(0x830c41c)])               |
2|  st  |>>Matched production: ['iEt' st]<<    |
2|  st  |>>Matched rule<< (return value:        |
|      |[ARRAY(0x830c41c)])               |
2|  st  |(consumed: [iEt iEt o e o])         |
1| prog |>>Matched subrule: [st]<< (return      |
|      |value: [ARRAY(0x830c41c)])         |

```

A estas alturas hemos construido el árbol:

```

st
|
' 'iEt' st
      |
      ' 'iEt' st      'e' st
            |          |
            ' 'o'      'o'

```

y el punto y coma nos espera en la entrada:

```

1| prog |Trying terminal: [';']              |
1| prog |>>Matched terminal<< (return value: |
|      |[;])                              |
1| prog |                                  |"\n"
1| prog |Trying action                      |
1| prog |>>Matched action<< (return value:   |
|      |[ARRAY(0x830c314)])               |
1| prog |>>Matched production: [st ';'<<    |
1| prog |>>Matched rule<< (return value:     |
|      |[ARRAY(0x830c314)])               |
1| prog |(consumed: [iEt iEt o e o ;])      |

```

Ahora ya se ejecuta la línea `print Dumper($result)` en el programa principal, volcando la estructura de datos construida durante el análisis:

```

$VAR1 = [ 'prog',
          [ 'st',

```

```

        'iEt',
        [ 'st',
          'iEt', 'o', 'e', 'o'
        ]
      ],
      ';'
    ];

```

Ejercicio 6.3.1. Si cambiamos el orden de las producciones y no forzamos a que la sentencia acabe en punto y coma:

```

my $grammar = q{
  st : 'iEt' st { [ @item ] }
    | 'iEt' st 'e' st { [ @item ] }
    | 'o'
};

```

¿Que árbol se obtendrá al darle la entrada `iEt iEt o e o ;`? ¿Que ocurre si, manteniendo este orden, forzamos a que el programa termine en punto y coma?

```

my $grammar = q{
  prog : st ';'
  st : 'iEt' st { [ @item ] }
    | 'iEt' st 'e' st { [ @item ] }
    | 'o'
};

```

¿Es aceptada la cadena?

Ejercicio 6.3.2. La regla “*if then else*” con el *if* mas cercano” puede hacerse explícita usando esta otra gramática:

<i>SENTENCIA</i>	→	<i>EQUI</i>		<i>NOEQUI</i>
<i>EQUI</i>	→	<i>if EXPRESION then EQUI else EQUI</i>		<i>OTRASSENTENCIAS</i>
<i>NOEQUI</i>	→	<i>if EXPRESION then SENTENCIA</i>		
		<i>if EXPRESION then EQUI else NOEQUI</i>		

Escriba un analizador sintáctico para la gramática usando *RD* y analice su comportamiento.

6.4. La directiva `commit`

El excesivo costo de *RD* puede aliviarse haciendo uso de las directivas `<commit>` y `<uncommit>`, las cuales permiten podar el árbol de búsqueda. Una directiva `<commit>` le indica al analizador que debe ignorar las producciones subsiguientes si la producción actual fracasa. Es posible cambiar esta conducta usando posteriormente la directiva `<uncommit>`, la cual revoca el estatus del último `<commit>`. Así la gramática del *if-then-else* puede reescribirse como sigue:

```

$ cat -n ifthenelse_commit.pl
1  #!/usr/bin/perl -w
2  use strict;
3  use Parse::RecDescent;
4  use Data::Dumper;
5
6  $::RD_TRACE = 1;
7  $::RD_AUTOACTION = q{ [@item] };

```

```

8 my $grammar = q{
9   prog : st ';'
10  st : 'iEt' <commit> st <uncommit> 'e' st
11      | 'iEt' <commit> st
12      | 'o' { 'o' }
13 };
14 ...

```

en este caso, si no se encuentra una sentencia después de `iEt` se producirá directamente el fracaso de `st` sin intentar las otras dos reglas. Sin embargo, si no se encuentra un `e` después de la sentencia si que se va a intentar la regla de la línea 11, ya que la directiva `<uncommit>` revoca el `<commit>` previo.

Observe el resultado de la ejecución: la acción automática de construcción del árbol da lugar a la inserción de los valores devueltos por las directivas `<commit>`:

```
./ifthenelse_commit.pl file4.txt
```

```
iEt iEt o e o ;
```

```

$VAR1 = [ 'prog',
           [ 'st',
             'iEt', 1, [ 'st',
                        'iEt', 1, 'o', 1, 'e', 'o'
                      ]
           ],
           ';'
         ];

```

Ejercicio 6.4.1. *¿Cuanto ahorro produce en este caso, el uso de las directivas `<commit>` en la gramática anterior? (supuesto que la entrada es `iEt iEt o e o ;`) ¿Se obtiene algún beneficio?*

6.5. Las Directivas `skip` y `leftop`

Consideremos una extensión al conocido problema del reconocimietno del lenguaje de *valores separados por comas* (véase la sección 3.4.3). El fichero de entrada contiene líneas como esta:

```
"tierra, luna",1,"marte, deimos",9.374
```

la extensión consiste en admitir que el separador puede cambiar entre líneas y en los campos de una línea, pudiendo ser no sólo la coma sino también un tabulador o el carácter dos puntos. El problema está en que los separadores dentro de las cadenas no deben confundirnos. Por ejemplo, la coma en `"tierra, luna"` no debe ser interpretada como un separador.

He aquí una solución usando RD:

```

1 #!/usr/local/bin/perl5.8.0 -w
2 use strict;
3 use Parse::RecDescent;
4 use Data::Dumper;
5
6 my $grammar = q{
7   line : <leftop: value sep value>
8   value: /"([^\"]*)" / { $1; }
9   sep  : <skip:""> /[:,\t]/
10 };
11

```



```

12 my $parse = Parse::RecDescent->new($grammar);
13
14 my $line;
15 while ($line = <>) {
16     print "$line\n";
17     my $result = $parse->line($line);
18     print Dumper($result);
19 }

```

La directiva `<skip>` usada en la línea 9 permite redefinir lo que el analizador entiende por “blancos”: los símbolos sobre los que saltar. Por defecto el tabulador es un “blanco”. En este ejemplo, cuando estamos reconociendo un separador no queremos que sea así. La directiva `<skip:"">` hace que ningún símbolo sea considerado como blanco.

La directiva usada en la línea 7 `<leftop: value sep value>` especifica una lista asociativa a izquierdas de elementos de tipo `value` separados por elementos del lenguaje denotado por `sep`. El valor retornado es un array anónimo conteniendo la lista de elementos. En general, si el separador es especificado como una cadena, la directiva no retorna el separador como parte de la lista. Por ejemplo, en el código:

```
list: '(' <leftop: list_item ',' list_item> ')' { $return = $item[2] }
```

el valor devuelto es una referencia a la lista, la cual no incluye las comas. Sin embargo, cuando el separador viene dado por una variable sintáctica o bien por medio de una expresión regular, las distintas apariciones del separador se incluyen en la lista devuelta. Este es el caso en nuestro ejemplo.

Veamos primero el fichero de entrada para la ejecución:

```

$ cat file.txt
"xx":"yy":"zzz"
"xx","yy","zzz"
"xx"    "yy"    "zzz"
"xx"    "yy","zzz"
"xx":"yy","zzz"
"x:x":"y,y","z  zz"
"x:x":"y,y","z  zz"."ttt"

```

La última fila usa el punto como separador del último campo. La ejecución de `Dumper($result)` en la línea 18 nos muestra las listas formadas:

```

$ ./commaseparated.pl file.txt
"xx":"yy":"zzz"
$VAR1 = [ 'xx', ':', 'yy', ':', 'zzz' ];
"xx","yy","zzz"
$VAR1 = [ 'xx', ',', 'yy', ',', 'zzz' ];
"xx"    "yy"    "zzz"
$VAR1 = [ 'xx', ' ', 'yy', ' ', 'zzz' ];
"xx"    "yy","zzz"
$VAR1 = [ 'xx', ' ', 'yy', ',,', 'zzz' ];
"xx":"yy","zzz"
$VAR1 = [ 'xx', ':', 'yy', ',,', 'zzz' ];
"x:x":"y,y","z  zz"
$VAR1 = [ 'x:x', ':', 'y,y', ',,', 'z  zz' ];
"x:x":"y,y","z  zz"."ttt"
$VAR1 = [ 'x:x', ':', 'y,y', ',,', 'z  zz' ];

```

La última entrada `"x:x":"y,y","z zz"."ttt"` muestra como el analizador se detiene en la cadena máxima `"x:x":"y,y","z zz"` que pertenece al lenguaje.

Ejercicio 6.5.1. *La última salida del programa anterior no produce un rechazo de la cadena de entrada. Se limita a detenerse en la cadena máxima "x:x":"y,y","z zz" que puede aceptar. Si queremos que la conducta sea de rechazo, se puede hacer que sea la línea completa la que se empareje. Escriba una nueva versión del programa anterior que recoja esta especificación.*

Ejercicio 6.5.2. *Escriba una variante del ejercicio anterior en la que se fuerze a que todos los separadores en una misma línea sean iguales (aunque pueden ser diferentes en diferentes líneas).*

6.6. Las directivas rulevar y reject

En los apuntes de Conway del “Advanced Perl Parsing” realizado durante el Deutscher Perl-Workshop 3.0 en el año 2001 se plantea la siguiente variante del problema de los valores separados por comas: los campos no están delimitados por dobles comillas y el separador es aquél que separa la línea en un mayor número de elementos, o lo que es lo mismo, aquél separador (de entre un grupo fijo `[:,\t]`) que se repite mas veces en la línea. Por ejemplo, la línea:

```
a:1,b:2,c:3,d
```

se separará por comas produciendo la lista `["a:1", "b:2","c:3","d"]`, pero la línea

```
a:1,b:2,c:3:d
```

se separará mediante el carácter dos puntos produciendo la lista `["a","1,b", "2,c","3","d"]`.

Para resolver el problema se declaran dos variables `max` y `$maxcount` locales al método asociado con la variable sintáctica `line` que reconoce la categoría gramatical línea. Para declarar tales variables se usa la directiva `rulevar`.

La otra directiva que se usa es `reject`, la cuál hace que la producción falle (exactamente igual que si la acción hubiese retornado `undef`). La estrategia consiste en almacenar en, separador por separador llevar en `$maxcount` el máximo número de items en que se descompone la línea y en `$max` la referencia al array anónimo “ganador”.

```
#!/usr/local/bin/perl5.8.0 -w
use strict;
use Parse::RecDescent;
use Data::Dumper;

#$::RD_TRACE = 1;
my $grammar = q{
    line : <rulevar: $max>
    line : <rulevar: $maxcount = 0>
    line : <leftop: value ',' value>
        { $maxcount = @{$item[1]}; $max = $item[1];
          print "maxcount[,] = $maxcount\n"; }
        <reject>
    | <leftop: datum ':' datum>
        { if (@{$item[1]} > $maxcount) {
            $maxcount = @{$item[1]}; $max = $item[1]; }
          print "maxcount[:,] = $maxcount\n";
        }
        <reject>
    | <leftop: field ";" field>
        { if (@{$item[1]} > $maxcount) {
            $maxcount = @{$item[1]}; $max = $item[1]; }
          print "maxcount[:,:] = $maxcount\n";
        }
}
```

```

        <reject>
        | { $return = $max; }

    value: /[^\,]*/
    datum: /[^\:]*/*
    field: /[^\;]*/*
};

my $parse = Parse::RecDescent->new($grammar);
my $line;
while ($line = <>) {
    print "$line\n";
    my $result = $parse->line($line);
    if (defined($result)) { print Dumper($result); }
    else { print "Cadena no válida\n"; }
}

```

Probemos el código con el siguiente fichero de entrada:

```

cat file3.txt
1:2:3,3:4;44
1,2:3,3:4;44
1;2:3;3:4;44

```

He aquí una ejecución:

```

$ ./maxseparator.pl file3.txt
1:2:3,3:4;44

maxcount[,] = 2
maxcount[:,] = 4
maxcount[:,:] = 4
$VAR1 = [ '1', '2', '3,3', '4;44' ];
1,2:3,3:4;44

maxcount[,] = 3
maxcount[:,] = 3
maxcount[:,:] = 3
$VAR1 = [ '1', '2:3', '3:4;44' ];
1;2:3;3:4;44

maxcount[,] = 1
maxcount[:,] = 3
maxcount[:,:] = 4
$VAR1 = [ '1', '2:3', '3:4', '44' ];

```

6.7. Utilizando score

La directiva `<score: ...>` toma como argumento una expresión numérica que indica la prioridad de la producción. La directiva lleva un `<reject>` asociado. La producción con puntuación mas alta es elegida (salvo que hubiera otra que hubiera aceptado directamente). El problema anterior puede ser resuelto como sigue:

```

#!/usr/local/bin/perl5.8.0 -w
use strict;

```

```

use Parse::RecDescent;
use Data::Dumper;

#$:RD_TRACE = 1;
my $grammar = q{
    line : <lefttop: value ',' value>
          <score: @{$item[1]}>
        | <lefttop: datum ':' datum>
          <score: @{$item[1]}>
        | <lefttop: field ";" field>
          <score: @{$item[1]}>

    value: /[^\,]*/
    datum: /^[^:]*/*
    field: /^[^;]*/*
};

my $parse = Parse::RecDescent->new($grammar);

my $line;
while ($line = <>) {
    print "$line\n";
    my $result = $parse->line($line);
    if (defined($result)) { print Dumper($result); }
    else { print "Cadena no válida\n"; }
}

```

Al darle como entrada el fichero

```

$ cat file3.txt
1:2:3,3:4;44
1,2:2,3:3;33
1;2:2;3:3;44

```

Obtenemos la ejecución:

```

$ ./score.pl file3.txt
1:2:3,3:4;44
$VAR1 = [ '1', '2', '3,3', '4;44 ' ];
1,2:2,3:3;33
$VAR1 = [ '1', '2:2', '3:3;33 ' ];
1;2:2;3:3;44
$VAR1 = [ '1', '2:2', '3:3', '44 ' ];

```

6.8. Usando autoscore

La solución anterior puede simplificarse usando **autoscore**. Si una directiva **<autoscore>** aparece en cualquier producción de una variable sintáctica, el código especificado por la misma es utilizado como código para la puntuación de la producción dentro de las producciones de la variable sintáctica. Con la excepción, claro está, de aquellas reglas que tuvieran ya una directiva **score** explícitamente asignada. El código de la gramática queda como sigue:

```

my $grammar = q{
    line : <autoscore: @{$item[1]}>
          | <lefttop: value ',' value>

```

```

        | <leftop: datum ':' datum>
        | <leftop: field ";" field>

value: /[^\,]*/
datum: /[^\: ]*/
field: /[^\;]*/
};

```

Por defecto, el analizador generado por `parse::RecDescent` siempre acepta la primera producción que reconoce un prefijo de la entrada. Como hemos visto, se puede cambiar esta conducta usando la directiva `<score: ...>`. Las diferentes producciones son intentadas y su puntuación (`score`) evaluada, considerándose vencedora aquella con mayor puntuación. Esto puede usarse para hacer que la regla vencedora sea “la mas larga”, en el sentido de ser la que mas elementos en `@item` deja. En este caso tomamos ventaja de la directiva `<autoscore>`, la cual permite asociar una directiva `<score>` con cada producción que no tenga ya una directiva `<score>` explícita.

```

#!/usr/local/bin/perl5.8.0 -w
use strict;
use warnings;
use Parse::RecDescent;

$::RD_TRACE = 1;
#$::RD_HINT = 1;
my $grammar = q {

    start:  seq_1 ';'
    seq_1   : <autoscore: @item>
              | 'A' 'B' 'C' 'D'
                { local $^W = 0;
                  print "seq_1: " . join (" ", @item[1..$#item]) . " score: $score\n"
                }
              | 'A' 'B' 'C' 'D' 'E' 'F'
                { print "seq_1: " . join (" ", @item[1..$#item]) . " score: $score\n" }
};

my $parser;

{
    local $^W = 0; # Avoid warning "uninitialized value in (m//) at RecDescent.pm line 626."
    $parser=Parse::RecDescent->new($grammar);
}

my $result = $parser->start("A B C D E F ;");
if (defined($result)) {
    print "Valida\n";
} else { print "No reconocida\n"; }

```

6.9. El Hash %item

En el contexto de una regla de producción, el hash `%item` proporciona un acceso indexado en los símbolos a los correspondientes atributos. Así, si tenemos la regla `A : B C D`, el elemento `$item{B}` hace alusión al atributo de B e `$item{C}` hace alusión al atributo de C. Los símbolos correspondientes a cadenas, expresiones regulares, directivas y acciones se guardan respectivamente bajo una clave de la forma `__STRINGn__`, `__PATTERNn__`, `__DIRECTIVEN__`, `__ACTIONn__`, donde `n` indica la posición

ordinal del elemento dentro de los de su tipo en la regla. Así, `__PATTERN2__` hace alusión a la segunda expresión regular en la parte derecha de la regla.

El elemento especial `$item{__RULE__}` contiene el nombre de la regla actual.

La ventaja de usar `%item` en vez de `@items` es que evita los errores cometidos al introducir o eliminar elementos en la parte derecha de una producción. Cuando esto ocurre, el programador debe cambiar los números de los índices en la acción.

Una limitación del hash `%item` es que, cuando hay varias apariciones de una variable sintáctica sólo guarda el valor de la última aparición. Por ejemplo:

```
range: '(' number '..' number ')' { $return = $item{number} }
```

retorna en `$item{number}` el valor correspondiente a la segunda aparición de la variable sintáctica `number`.

6.10. Usando la directiva autotree

La directiva `<autotree>` hace que RD construya automáticamente una representación del árbol sintáctico concreto. Cada nodo del árbol, correspondiendo a una parte derecha de una regla de producción es bendecido (usando `bless`) en una clase cuyo nombre es el de la variable sintáctica que produce: viene a ser equivalente a `{ bless \%item, $item[0] }`.

```
$ cat ./infixautotree.pl
#!/usr/local/bin/perl5.8.0 -w
use strict;
use Parse::RecDescent;
use Data::Dumper;

#$:RD_TRACE = 1;
my $grammar = q{
    <autotree>
    expre  : <leftop: prods '+' prods>
    prods  : <leftop: unario '*' unario>
    unario: "(" expre ")" | numero
    numero : /\d+(\.\d+)?/
};

my $parse = new Parse::RecDescent($grammar);
my $line = shift;
my $result = $parse->expre($line);
if (defined($result)) {
    print "Válida:\n", Dumper($result);
}
else { print "Cadena no válida\n"; }
```

Las reglas que consisten en un sólo terminal, como la de número, llevan un tratamiento especial. En ese caso el código ejecutado es:

```
{ bless {__VALUE__=>$item[1]}, $item[0] }
```

La estructura resultante es relativamente compleja. Cada nodo es bendecido en la clase con nombre el de la variable sintáctica, esto es, con nombre el de la regla. Podemos sacar ventaja de esta aproximación si asociamos métodos de procesamiento con cada tipo de nodo. Por ejemplo, para hacer un recorrido del árbol sintáctico, podríamos extender el ejemplo anterior como sigue:

```

$ cat infixautotree2.pl
...

sub expre::traverse {
    my $root = shift;
    my $level = shift;

    process($root, $level, '__RULE__');
    foreach (@{$root->{__DIRECTIVE1__}}) {
        $_->traverse($level+1);
    }
}

sub prods::traverse {
    my $root = shift;
    my $level = shift;

    process($root, $level, '__RULE__');
    foreach (@{$root->{__DIRECTIVE1__}}) {
        $_->traverse($level+1);
    }
}

sub unario::traverse {
    my $root = shift;
    my $level = shift;

    process($root, $level, '__RULE__');
    my $child = $root->{numero} || $root->{expre};
    $child->traverse($level+1);
}

sub numero::traverse {
    my $root = shift;
    my $level = shift;

    process($root, $level, '__VALUE__');
}

sub process {
    my $root = shift;
    my $level = shift;
    my $value = shift;

    print " "x($width*$level),$root->{$value},"\\n";
}

```

Sigue un ejemplo de ejecución:

```
$ ./infixautotree2.pl '2+3*(5+2)+9' 6
```

Válida:

expre

 prods

 unario

```

      2
    prods
      unario
        3
      unario
        expre
          prods
            unario
              5
            prods
              unario
                2
          prods
            unario
              9

```

6.11. Práctica

Reescriba la fase de análisis sintáctico del compilador para el lenguaje Tutu usando `Parse::RecDescent`. La gramática de Tutu es como sigue:

```

program → block
block → declarations statements | statements
declarations → declaration ';' declarations | declaration ';'
declaration → INT idlist | STRING idlist
statements → statement ';' statements | statement
statement → ID '=' expression | P expression | '{' block '}'
expression → expression '+' term | expression '-' term | term
term → term '*' factor | term '/' factor | factor
factor → '(' expression ')' | ID | NUM | STR
idlist → ID ',' idlist | ID

```

Las acciones deberán producir el árbol de análisis abstracto o AST. Recicle todo el código que ha escrito para las restantes fases: análisis léxico, semántico, optimización de código, cálculo de direcciones, generación de código y optimización peephole.

6.12. Construyendo un compilador para Parrot

Las ideas y el código en esta sección están tomadas del artículo de Dan Sugalski *Building a parrot Compiler* que puede encontrarse en <http://www.onlamp.com/lpt/a/4725>.

```

1 #!/usr/local/bin/perl5.8.0 -w
2 #
3 # This program created 2004, Dan Sugalski. The code in this file is in
4 # the public domain--go for it, good luck, don't forget to write.
5 use strict;
6 use Parse::RecDescent;
7 use Data::Dumper;
8
9 # Take the source and destination files as parameters
10 my ($source, $destination) = @ARGV;
11
12 my %global_vars;
13 my $tempcount = 0;

```



```

14 my (%temps) = (P => 0,
15               I => 0,
16               N => 0,
17               S => 0
18               );
19
20 # AUTOACTION simplifies the creation of a parse tree by specifying an action
21 # for each production (ie action is { [@item] })
22 $::RD_AUTOACTION = q{ [@item] };
23
24 my $grammar = <<'EOG';
25 field: /\b\w+\b/
26
27 stringconstant: /'[^']*'/ |
28               /"[^"]*" /
29 #"
30 float: /[+-]?(?=\d|\.\d)\d*(\.\d*)?([Ee]([+-]?\d+))?/
31
32 constant: float | stringconstant
33
34 addop: '+' | '-'
35 mulop: '*' | '/'
36 modop: '%'
37 cmpop: '<>' | '>=' | '<=' | '<' | '>' | '='
38 logop: 'and' | 'or'
39
40 parenexpr: '(' expr ')'
41
42 simplevalue: parenexpr | constant | field
43
44 modval: <leftop: simplevalue modop simplevalue>
45
46 mulval: <leftop: modval mulop modval>
47
48 addval: <leftop: mulval addop mulval>
49
50 cmpval: <leftop: addval cmpop addval>
51
52 logval: <leftop: cmpval logop cmpval>
53
54 expr: logval
55
56 declare: 'declare' field
57
58 assign: field '=' expr
59
60 print: 'print' expr
61
62 statement: assign | print | declare
63 EOG
64
65 # ?? Makes emacs cperl syntax highlighting mode happier
66 my $parser = Parse::RecDescent->new($grammar);

```

La gramática categoriza las prioridades de cada una de las operaciones: categorías próximas al símbolo de arranque tienen menos prioridad que aquellas más lejanas.

```

68 my @nodes;
69 open SOURCE, "<$source" or die "Can't open source program ($!)";
70
71 while (<SOURCE>) {
72     # Strip the trailing newline and leading spaces. If the line is
73     # blank, then skip it
74     chomp;
75     s/^\s+//;
76     next unless $_;
77
78     # Parse the statement and throw an error if something went wrong
79     my $node = $parser->statement($_);
80     die "Bad statement" if !defined $node;
81
82     # put the parsed statement onto our list of nodes for later treatment
83     push @nodes, $node;
84 }
85
86 print Dumper(\@nodes);
87 #exit;
88
89 # At this point we have parsed the program and have a tree of it
90 # ready to process. So let's do so. First we set up our node handlers.
91

```

El programa principal lee una línea del fuente y construye el árbol (línea 79). Los árboles se van guardando en la lista @nodes. El paso siguiente es la generación de código:

```

# At this point we have parsed the program and have a tree of it
# ready to process. So let's do so. First we set up our node handlers.

```

```

my (%handlers) = (addval => \&handle_generic_val,
    assign => \&handle_assign,
    cmpval => \&handle_generic_val,
    constant => \&delegate,
    declare => \&handle_declare,
    expr => \&delegate,
    field => \&handle_field,
    float => \&handle_float,
    logval => \&handle_generic_val,
    modval => \&handle_generic_val,
    mulval => \&handle_generic_val,
    negfield => \&handle_negfield,
    parenexpr => \&handle_paren_expr,
    print => \&handle_print,
    simplevalue => \&delegate,
    statement => \&delegate,
    stringconstant => \&handle_stringconstant,
);

```

El autor ha optado por escribir un manipulador para cada tipo de nodo. Es algo similar a lo que hicimos usando métodos y herencia para el compilador de Tutu. Algunos nodos simplemente delegan y otros recurren a un manipulador genérico.

La fase de generación de código comienza por la escritura de un preámbulo y termina con la escritura de un pie requeridos por el intérprete. En medio se sitúa el código correspondiente a la traducción de los nodos provenientes de las diversas líneas del fuente:

```
# Open the output file and emit the preamble
open PIR, ">$destination" or die "Can't open destination ($!)";
print PIR <<HEADER;
.sub __MAIN prototyped
    .param pmc argv
HEADER

foreach my $node (@nodes) {
    my @lines = process_node(@$node);
    print PIR join("", @lines);
}

print PIR <<FOOTER;
    end
.end
FOOTER
```

La subrutina `process_node` hace un recorrido de los árboles de análisis, llamando a los manipuladores de los nodos que están siendo visitados. El elemento 0 del array `elems` identifica la clase de nodo. Así la llamada `$handlers{$elems[0]}->(@elems)` produce una llamada al manipulador correspondiente, pasándole como argumento los hijos del nodo.

```
# The value of the last expression evaluated
sub last_expr_val {
    return $::last_expr;
}

# Setting the last expression evaluated's value
sub set_last_expr_val {
    $::last_expr = $_[0];
}

sub process_node {
    my (@elems) = @_;
    return "\n" unless @elems;
    return "\n" unless defined($elems[0]);
    if (ref $elems[0]) {
        return process_node(@{$elems[0]});
    } elsif (exists($handlers{$elems[0]})) {
        return $handlers{$elems[0]}->(@elems);
    } else {
        return "****", $elems[0], "****\n";
    }
}
```

A continuación siguen los diversos manipuladores para los diferentes tipos de nodo:

```
sub handle_assign {
    my ($nodetype, $destvar, undef, $expr) = @_;
    my @nodes;
    push @nodes, process_node(@$expr);
}
```

```

    my $rhs = last_expr_val();
    push @nodes, process_node(@$destvar);
    my $lhs = last_expr_val();
    push @nodes, "    $lhs = $rhs\n";
    return @nodes;
}

sub handle_declare {
    my ($nodetype, undef, $var) = @_;
    my @lines;

    my $varname = $var->[1];

    # Does it exist?
    if (defined $global_vars{$varname}) {
        die "Multiple declaration of $varname";
    }
    $global_vars{$varname}++;
    push @lines, "    .local pmc $varname\n";
    push @lines, "    new $varname, .PerlInt\n";
    return @lines;
}

sub handle_field {
    my ($nodetype, $fieldname) = @_;
    if (!exists $global_vars{$fieldname}) {
        die "undeclared field $fieldname used";
    }
    set_last_expr_val($fieldname);
    return;
}

sub handle_float {
    my ($nodetype, $floatval) = @_;
    set_last_expr_val($floatval);
    return;
}

sub handle_generic_val {
    my (undef, $terms) = @_;
    my (@terms) = @$terms;

    # Process the LHS
    my $lhs = shift @terms;
    my @tokens;
    push @tokens, process_node(@$lhs);

    my ($op, $rhs);

    # Now keep processing the RHS as long as we have it
    while (@terms) {
        $op = shift @terms;
        $rhs = shift @terms;
    }
}

```

```

my $val = last_expr_val();
my $oper = $op->[1];

push @tokens, process_node(@$rhs);
my $other_val = last_expr_val();

my $dest = $temps{P}++;

foreach ($oper) {
    # Simple stuff -- addition, subtraction, multiplication,
    # division, and modulus. Just a quick imcc transform
    /(\+|\-|\*|\/|\%)/ && do { push @tokens, "new \$$dest, .PerlInt\n";
        push @tokens, "\$$dest = $val $oper $other_val\n";
        set_last_expr_val("\$$dest");
        last;
    };
    /and/ && do { push @tokens, "new \$$dest, .PerlInt\n";
        push @tokens, "\$$dest = $val && $other_val\n";
        set_last_expr_val("\$$dest");
        last;
    };
    /or/ && do { push @tokens, "new \$$dest, .PerlInt\n";
        push @tokens, "\$$dest = $val || $other_val\n";
        set_last_expr_val("\$$dest");
        last;
    };
    /</ && do { my $label = "eqcheck$tempcount";
        $tempcount++;
        push @tokens, "new \$$dest, .Integer\n";
        push @tokens, "\$$dest = 1\n";
        push @tokens, "ne $val, $other_val, $label\n";
        push @tokens, "\$$dest = 0\n";
        push @tokens, "$label:\n";
        set_last_expr_val("\$$dest");
        last;
    };
    /=/ && do { my $label = "eqcheck$tempcount";
        $tempcount++;
        push @tokens, "new \$$dest, .Integer\n";
        push @tokens, "\$$dest = 1\n";
        push @tokens, "eq $val, $other_val, $label\n";
        push @tokens, "\$$dest = 0\n";
        push @tokens, "$label:\n";
        set_last_expr_val("\$$dest");
        last;
    };
    />/ && do { my $label = "eqcheck$tempcount";
        $tempcount++;
        push @tokens, "new \$$dest, .Integer\n";
        push @tokens, "\$$dest = 1\n";
        push @tokens, "lt $val, $other_val, $label\n";
        push @tokens, "\$$dest = 0\n";
        push @tokens, "$label:\n";
    };
}

```

```

        set_last_expr_val("\$P$dest");
        last;
    };
    />/ && do { my $label = "eqcheck$tempcount";
        $tempcount++;
        push @tokens, "new \$P$dest, .Integer\n";
        push @tokens, "\$P$dest = 1\n";
        push @tokens, "gt $val, $other_val, $label\n";
        push @tokens, "\$P$dest = 0\n";
        push @tokens, "$label:\n";
        set_last_expr_val("\$P$dest");
        last;
    };
    die "Can't handle $oper";
}

}
return @tokens;
}

sub handle_paren_expr {
    my ($nodetype, undef, $expr, undef) = @_;
    return process_node(@$expr);
}

sub handle_stringconstant {
    my ($nodetype, $stringval) = @_;
    set_last_expr_val($stringval);
    return;
}

sub handle_print {
    my ($nodetype, undef, $expr) = @_;
    my @nodes;
    push @nodes, process_node(@$expr);
    my $val = last_expr_val();
    push @nodes, " print $val\n";
    return @nodes;
}

sub delegate {
    my ($nodetype, $nodeval) = @_;
    return process_node(@$nodeval);
}

```

El fichero fuente foo.len:

```

declare foo
declare bar

foo = 15
bar = (foo+8)*32-7

print bar
print "\n"

```

```
print foo % 10
print "\n"
```

Compilamos:

```
$ ./compiler.pl foo.len foo.pir
```

Esto produce por pantalla un volcado de los árboles de las diferentes sentencias. Así para `declare foo`:

```
$VAR1 = [ [ 'statement', [ 'declare', 'declare', [ 'field', 'foo' ] ] ],
```

para la sentencia `foo = 15` el árbol es:

```
[ 'statement',  
  [ 'assign',  
    [ 'field',  
      'foo'  
    ],  
    '=',  
    [ 'expr',  
      [ 'logval',  
        [ 'cmpval',  
          [ 'addval',  
            [ 'mulval',  
              [ 'modval',  
                [ 'simplevalue',  
                  [ 'constant',  
                    [ 'float',  
                      '15'  
                    ]  
                  ]  
                ]  
              ]  
            ]  
          ]  
        ]  
      ]  
    ]  
  ] ] ] ] ] ] ] ] ] ] ] ] ],
```

Este es el árbol de la sentencia `print bar`:

```
[ 'statement',
  [ 'print',
    'print',
    [ 'expr',
      [ 'logval',
        [
          [ 'cmpval',
            [
              [ 'addval',
                [
                  [ 'mulval',
                    [
                      [ 'modval',
                        [
```



```

    --gc-debug
    --leak-test|--destroy-at-end
    -. --wait    Read a keystroke before starting
<Compiler options>
    -v --verbose
    -E --pre-process-only
    -o --output=FILE
        --output-pbc
    -O --optimize[=LEVEL]
    -a --pasm
    -c --pbc
    -r --run-pbc
    -y --yydebug
<Language options>
    --python

```

see docs/running.pod for more

Con la opción `-o` podemos producir un fichero en formato `pbc`:

```
$ parrot -o foo.pbc foo.pir
```

que podemos ejecutar con el depurador `pdb` (para construirlo en el momento de la instalación de Parrot deberás hacer `make pdb`).

```
$ pdb foo.pbc
```

Parrot Debugger 0.0.2

```
(pdb) h
```

List of commands:

```

disassemble -- disassemble the bytecode
load -- load a source code file
list (l) -- list the source code file
run (r) -- run the program
break (b) -- add a breakpoint
watch (w) -- add a watchpoint
delete (d) -- delete a breakpoint
disable -- disable a breakpoint
enable -- reenable a disabled breakpoint
continue (c) -- continue the program execution
next (n) -- run the next instruction
eval (e) -- run an instruction
trace (t) -- trace the next instruction
print (p) -- print the interpreter registers
stack (s) -- examine the stack
info -- print interpreter information
quit (q) -- exit the debugger
help (h) -- print this help

```

Type "help" followed by a command name for full documentation.

Veamos el programa traducido:

```
(pdb) list 1 17
```

```

1  new_p_ic P16,32
2  new_p_ic P30,32
3  set_p_ic P16,15

```

```

4  new_p_ic P29,32
5  add_p_p_ic P29,P16,8
6  new_p_ic P28,32
7  mul_p_p_ic P28,P29,32
8  new_p_ic P29,32
9  sub_p_p_ic P29,P28,7
10 set_p_p P30,P29
11 print_p P30
12 print_sc "\n"
13 new_p_ic P30,32
14 mod_p_p_ic P30,P16,10
15 print_p P30
16 print_sc "\n"
17 end

```

Procedemos a ejecutarlo:

```

(pdb) n
2  new_p_ic P30,32
(pdb)
3  set_p_ic P16,15
(pdb)
4  new_p_ic P29,32
(pdb)
5  add_p_p_ic P29,P16,8
(pdb) p P16
P16 = [PerlInt]
Stringified: 15
5  add_p_p_ic P29,P16,8
(pdb) c
729
5
Program exited.
(pdb) quit
$

```

6.13. Práctica

Extienda el lenguaje Tutu para que comprenda el tipo lógico, declarado mediante la palabra reservada `bool` y operaciones de comparación como las que figuran en el ejemplo propuesto por Sugalski. Genere código para la máquina basada en registros. Extienda el juego de instrucciones de manera apropiada.

6.14. Práctica

Modifique la fase de generación de código del compilador de Tutu para que produzca código para la máquina Parrot. Modifique el optimizador peephole para que se adapte al código Parrot.

Capítulo 7

Análisis LR

7.1. Parse::Yapp: Ejemplo de Uso

El generador de analizadores sintácticos `Parse::Yapp` es un analizador LALR inspirado en `yacc`. El módulo `Parse::Yapp` no viene con la distribución de Perl, por lo que es necesario bajarlo desde CPAN, en la dirección

`\http://search.cpan.org/~fdesar/Parse-Yapp-1.05/lib/Parse/Yapp.p`

o bien en nuestros servidores locales, por ejemplo en el mismo directorio en que se guarda la versión HTML de estos apuntes encontrará una copia de `Parse-Yapp-1.05.tar.gz`. La versión a la que se refiere este capítulo es la 1.05.

Para ilustrar su uso veamos un ejemplo en el que se genera una sencilla calculadora numérica. Los contenidos del programa `yapp` los hemos guardado en un fichero denominado `Calc.yp` (el código completo figura en el apéndice en la página ??)

```
1  #
2  # Calc.yp
3  #
4  # Parse::Yapp input grammar example.
5  #
6  # This file is PUBLIC DOMAIN
7  #
8  #
```

Se pueden poner comentarios tipo Perl o tipo C (`/* ... */`) a lo largo del fichero.

```
9  %right  '='
10 %left  '-' '+'
11 %left  '*' '/'
12 %left  NEG
13 %right  '^'
14
15 %%
16 input:  #empty
17         |  input line { push(@{$_[1]},$_[2]); $_[1] }
18         ;
```

las declaraciones `%left` y `%right` expresan la asociatividad y precedencia de los terminales, permitiendo decidir que árbol construir en caso de ambigüedad. *Los terminales declarados en líneas posteriores tienen mas prioridad que los declarados en las líneas anteriores.* Véase la sección 8.26 para mas detalles.

Un programa **yapp** consta de tres partes: la cabeza, el cuerpo y la cola. Cada una de las partes va separada de las otras por el símbolo **%%** en una línea aparte. Así, el **%%** de la línea 15 separa la cabeza del cuerpo. En la cabecera se colocan el código de inicialización, las declaraciones de terminales, las reglas de precedencia, etc. El cuerpo contiene las reglas de la gramática y las acciones asociadas. Por último, la cola de un program **yapp** contiene las rutinas de soporte al código que aparece en las acciones así como, posiblemente, rutinas para el análisis léxico y el tratamiento de errores.

En **Parse::Yapp** las acciones son convertidas en subrutinas anónimas. Mas bien en métodos anónimos. Así pues el primer argumento de la subrutina se identifica con una referencia al analizador (**\$_[0]**). Los restantes parámetros se corresponden con los *atributos de los símbolos* en la parte derecha de la regla de producción (**\$_[1] ...**). Por ejemplo, el código en la línea 21 imprime el atributo asociado con la variable sintáctica **expr**, que en este caso es su valor numérico.

La línea 17 indica que el atributo asociado con la variable sintáctica **input** es una referencia a una pila y que el atributo asociado con la variable sintáctica **line** debe empujarse en la pila. De hecho, el atributo asociado con **line** es el valor de la expresión. Así pues el atributo retornado por **input** es una referencia a una lista conteniendo los valores de las expresiones evaluadas.

Para saber mas sobre las estructuras internas de **yapp** para la representación de las acciones asociadas con las reglas véase la sección 7.4.

```

19
20 line:  '\n'          { 0 }
21      | exp '\n'      { print "$_[1]\n"; $_[1] }
22      | error '\n'    { $_[0]->YYErrork; 0 }
23 ;

```

El terminal **error** en la línea 22 esta asociado con la aparición de un error. El tratamiento es el mismo que en **yacc**. Cuando se produce un error en el análisis, **yapp** emite un mensaje de error y produce “mágicamente” un terminal especial denominado **error**. A partir de ahí permanecerá silencioso, consumiendo terminales hasta encontrar uno de los terminales que le hemos indicado en las reglas de recuperación de errores, en este caso, cuando encuentre un retorno de carro. Como se ha dicho, en **Parse::Yapp** el primer argumento de la acción denota al analizador sintáctico. Así pues el código **\$_[0]->YYErrork** es una llamada al método **YYErrork** del analizador. Este método funciona como la macro **yyerrork** de **yacc**, indicando que la presencia del retorno del carro (**\n**) la podemos considerar un signo seguro de que nos hemos recuperado del error. A partir de este momento, **yapp** volverá a emitir mensajes de error. Para saber más sobre la recuperación de errores en **yapp** léase la sección 7.15.

```

24
25 exp:      NUM

```

La acción por defecto es retornar **\$_[1]**. Por tanto, en este caso el valor retornado es el asociado a **NUM**.

```

26      | VAR          { $_[0]->YYData->{VARS}{$_[1]} }

```

El método **YYData** provee acceso a un hash que contiene los datos que están siendo analizados. En este caso creamos una entrada **VARS** que es una referencia a un hash en el que guardamos las variables. Este hash es la tabla de símbolos de la calculadora.

```

27      | VAR '=' exp   { $_[0]->YYData->{VARS}{$_[1]}=$_[3] }
28      | exp '+' exp   { $_[1] + $_[3] }
29      | exp '-' exp   { $_[1] - $_[3] }
30      | exp '*' exp   { $_[1] * $_[3] }

```

Hay numerosas ambigüedades en esta gramática. Por ejemplo,

- ¿Como debo interpretar la expresión $4 - 5 - 2$? ¿Como $(4 - 5) - 2$? ¿o bien $4 - (5 - 2)$? La respuesta la da la asignación de asociatividad a los operadores que hicimos en las líneas 9-13. Al declarar como asociativo a izquierdas al terminal $-$ (línea 10) hemos resuelto este tipo de ambigüedad. Lo que estamos haciendo es indicarle al analizador que a la hora de elegir entre los árboles abstractos $-(4, 5), 2$ y $-(4, -(5, 2))$ elija siempre el árbol que se hunde a izquierdas.
- ¿Como debo interpretar la expresión $4 - 5 * 2$? ¿Como $(4 - 5) * 2$? ¿o bien $4 - (5 * 2)$? Al declarar que $*$ tiene mayor prioridad que $-$ estamos resolviendo esta otra fuente de ambigüedad. Esto es así pues $*$ fué declarado en la línea 11 y $-$ en la 10.

```

31      |   exp '/' exp   {
32                               $_[3]
33      and return($_[1] / $_[3]);
34      $_[0]->YYData->{ERRMSG}
35      =   "Illegal division by zero.\n";
36      $_[0]->YYError;
37      undef
38      }
```

En la regla de la división comprobamos que el divisor es distinto de cero. Si es cero inicializamos el atributo `ERRMSG` en la zona de datos con el mensaje de error apropiado. Este mensaje es aprovechado por la subrutina de tratamiento de errores (véase la subrutina `_Error` en la zona de la cola). La subrutina `_Error` es llamada automáticamente por `yapp` cada vez que ocurre un error sintáctico. Esto es así por que en la llamada al analizador se especifican quienes son las diferentes rutinas de apoyo:

```

my $result = $self->YYParse( yylex => \&_Lexer,
                             yyerror => \&_Error,
                             yydebug => 0x0 );
```

Por defecto, *una regla de producción tiene la prioridad del último terminal que aparece en su parte derecha*. Una regla de producción puede ir seguida de una directiva `%prec` la cual le da una prioridad explícita. Esto puede ser de gran ayuda en ciertos casos de ambigüedad.

```

39      |   '-' exp %prec NEG   { -$_[2] }
```

¿Cual es la ambigüedad que surge con esta regla? Una de las ambigüedad de esta regla esta relacionada con el doble significado del menos como operador unario y binario: hay frases como $-y-z$ que tiene dos posibles interpretaciones: Podemos verla como $(-y)-z$ o bien como $-(y-z)$. Hay dos árboles posibles. El analizador, cuando este analizando la entrada $-y-z$ y vea el segundo $-$ deberá escoger uno de los dos árboles. ¿Cuál?. El conflicto puede verse como una “lucha” entre la regla `exp: '-' exp` la cual interpreta la frase como $(-y)-z$ y la segunda aparición del terminal $-$ el cuál “quiere entrar” para que gane la regla `exp: exp '-' exp` y dar lugar a la interpretación $-(y-z)$. En este caso, las dos reglas $E \rightarrow -E$ y $E \rightarrow E - E$ tienen, en principio la prioridad del terminal $-$, el cual fué declarado en la línea 10. La prioridad expresada explícitamente para la regla por la declaración `%prec NEG` de la línea 39 hace que la regla tenga la prioridad del terminal `NEG` (línea 12) y por tanto mas prioridad que el terminal $-$. Esto hará que `yapp` finalmente opte por la regla `exp: '-' exp`.

La declaración de \wedge como asociativo a derechas y con un nivel de prioridad alto resuelve las ambigüedades relacionadas con este operador:

```

40      |   exp '^' exp      { $_[1] ** $_[3] }
41      |   '(' exp ')'      { $_[2] }
42      ;
```

Después de la parte de la gramática, y separada de la anterior por el símbolo `%%`, sigue la parte en la que se suelen poner las rutinas de apoyo. Hay al menos dos rutinas de apoyo que el analizador

sintáctico requiere le sean pasados como argumentos: la de manejo de errores y la de análisis léxico. El método `Run` ilustra como se hace la llamada al método de análisis sintáctico generado, utilizando la técnica de llamada con argumentos con nombre y pasándole las referencias a las dos subrutinas (en Perl, es un convenio que si el nombre de una subrutina comienza por un guión bajo es que el autor la considera privada):

...

```
sub Run {
    my($self)=shift;
    my $result = $self->YYParse( yylex => \&_Lexer,
                                yyerror => \&_Error,
                                yydebug => 0x0 );

    my @result = @$result;
    print "@result\n";
}
```

La subrutina de manejo de errores `_Error` imprime el mensaje de error proveído por el usuario, el cual, si existe, fué guardado en `$_[0]->YYData->{ERRMSG}`.

```
43 # rutinas de apoyo
44 %%
45
46 sub _Error {
47     exists $_[0]->YYData->{ERRMSG}
48     and do {
49         print $_[0]->YYData->{ERRMSG};
50         delete $_[0]->YYData->{ERRMSG};
51         return;
52     };
53     print "Syntax error.\n";
54 }
55
```

A continuación sigue el método que implanta el análisis léxico `_Lexer`. En primer lugar se comprueba la existencia de datos en `parser->YYData->{INPUT}`. Si no es el caso, los datos se tomarán de la entrada estándar:

```
56 sub _Lexer {
57     my($parser)=shift;
58
59     defined($parser->YYData->{INPUT})
60     or $parser->YYData->{INPUT} = <STDIN>
61     or return('','undef');
62
```

Cuando el analizador léxico alcanza el final de la entrada debe devolver la pareja `('','undef')`. Eliminamos los blancos iniciales (lo que en inglés se conoce por *trimming*):

```
63 $parser->YYData->{INPUT} =~ s/^[ \t]//;
64
```

A continuación vamos detectando los números, identificadores y los símbolos individuales. El bucle `for ($parser->YYData->{INPUT})` se ejecuta mientras la cadena en `$parser->YYData->{INPUT}` no sea vacía, lo que ocurrirá cuando todos los terminales hayan sido consumidos.

```

65   for ($parser->YYData->{INPUT}) {
66       s/^[0-9]+(?:\.[0-9]+)?//
67       and return('NUM',$1);
68       s/^[A-Za-z][A-Za-z0-9_]*//
69       and return('VAR',$1);
70       s/^(.)//s
71       and return($1,$1);
72   }
73 }

```

Ejercicio 7.1.1. 1. ¿Quién es la variable índice en la línea 65?

2. ¿Sobre quién ocurre el binding en las líneas 66, 68 y 70?

3. ¿Cuál es la razón por la que `$parser->YYData->{INPUT}` se ve modificado si no aparece como variable para el binding en las líneas 66, 68 y 70?

Construimos el módulo `Calc.pm` a partir del fichero `Calc.y` especificando la gramática, usando un fichero `Makefile`:

```

> cat Makefile
Calc.pm: Calc.y
        yapp -m Calc Calc.y
> make
yapp -m Calc Calc.y

```

Esta compilación genera el fichero `Calc.pm` conteniendo el analizador:

```

> ls -ltr
total 96
-rw-r-----  1 pl      users      1959 Oct 20  1999 Calc.y
-rw-r-----  1 pl      users      39 Nov 16  12:26 Makefile
-rwxrwx--x   1 pl      users      78 Nov 16  12:30 usecalc.pl
-rw-rw----   1 pl      users     5254 Nov 16  12:35 Calc.pm

```

El script `yapp` es un *frontend* al módulo `Parse::Yapp`. Admite diversas formas de uso:

- `yapp [options] grammar[.yp]`

El sufijo `.yp` es opcional.

- `yapp -V`

Nos muestra la versión:

```

$ yapp -V
This is Parse::Yapp version 1.05.

```

- `yapp -h`

Nos muestra la ayuda:

```

$ yapp -h

```

```

Usage: yapp [options] grammar[.yp]
or     yapp -V
or     yapp -h

```

`-m module` Give your parser module the name `<module>`

```

                                default is <grammar>
-v                               Create a file <grammar>.output describing your parser
-s                               Create a standalone module in which the driver is included
-n                               Disable source file line numbering embedded in your parser
-o outfile                       Create the file <outfile> for your parser module
                                Default is <grammar>.pm or, if -m A::Module::Name is
                                specified, Name.pm
-t filename                      Uses the file <filename> as a template for creating the parser
                                module file. Default is to use internal template defined
                                in Parse::Yapp::Output
-b shebang                      Adds '#!<shebang>' as the very first line of the output file

grammar                          The grammar file. If no suffix is given, and the file
                                does not exists, .yp is added

-V                               Display current version of Parse::Yapp and gracefully exits
-h                               Display this help screen

```

La opción *-m module* da el nombre al paquete o espacio de nombres o clase encapsulando el analizador. Por defecto toma el nombre de la gramática. En el ejemplo podría haberse omitido.

La opción *-o outfile* da el nombre del fichero de salida. Por defecto toma el nombre de la gramática, seguido del sufijo *.pm*. sin embargo, si hemos especificado la opción *-m A::Module::Name* el valor por defecto será *Name.pm*.

Veamos los contenidos del ejecutable *usecalc.pl* el cuál utiliza el módulo generado por *yapp*:

```

> cat usecalc.pl
#!/usr/local/bin/perl5.8.0 -w

```

```
use Calc;
```

```

$parser = new Calc();
$parser->Run;

```

Al ejecutar obtenemos:

```

$ ./usecalc3.pl
2+3
5
4*8
32
^D
5 32

```

Pulsamos al final **Ctrl-D** para generar el final de fichero. El analizador devuelve la lista de valores computados la cual es finalmente impresa.

¿En que orden ejecuta **YYParse** las acciones? La respuesta es que el analizador generado por **yapp** construye una derivación a derechas inversa y ejecuta las acciones asociadas a las reglas de producción que se han aplicado. Así, para la frase **3+2** la antiderivación es:

$$NUM + NUM \xleftarrow{NUM \leftarrow E} E + NUM \xleftarrow{NUM \leftarrow E} E + E \xleftarrow{E + E \leftarrow E} E$$

por tanto las acciones ejecutadas son las asociadas con las correspondientes reglas de producción:

1. La acción de la línea 25:

25 `exp:` `NUM { $_[1]; } # acción por defecto`

Esta instancia de `exp` tiene ahora como atributo 3.

2. De nuevo la acción de la línea 25:

25 `exp:` `NUM { $_[1]; } # acción por defecto`

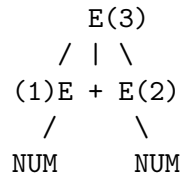
Esta nueva instancia de `exp` tiene como atributo 2.

3. La acción asociada con $E \rightarrow E + E$, en la línea 28:

28 `| exp '+' exp { $_[1] + $_[3] }`

La nueva instancia (nodo) `exp` tiene como atributo $5 = 3 + 2$.

Obsérvese que la antiderivación a derechas da lugar a un recorrido ascendente y a izquierdas del árbol:



Los números entre paréntesis indican el orden de visita de las producciones.

7.2. Conceptos Básicos

Los analizadores generador por `yapp` entran en la categoría de analizadores *LR*. Estos analizadores construyen una derivación a derechas inversa (o *antiderivación*). De ahí la *R* en *LR* (del inglés *rightmost derivation*). El árbol sintáctico es construido de las hojas hacia la raíz, siendo el último paso en la antiderivación la construcción de la primera derivación desde el símbolo de arranque.

Empezaremos entonces considerando las frases que pueden aparecer en una derivación a derechas. Tales frases consituyen el lenguaje *FSD*:

Definición 7.2.1. Dada una gramática $G = (\Sigma, V, P, S)$ no ambigua, se denota por *FSD* (lenguaje de las formas Sentenciales a Derechas) al lenguaje de las sentencias que aparecen en una derivación a derechas desde el símbolo de arranque.

$$FSD = \left\{ \alpha \in (\Sigma \cup V)^* : \exists S \xRightarrow[RM]{*} \alpha \right\}$$

Donde la notación *RM* indica una derivación a derechas (*rightmost*). Los elementos de *FSD* se llaman “formas sentenciales derechas”.

Dada una gramática no ambigua $G = (\Sigma, V, P, S)$ y una frase $x \in L(G)$ el proceso de antiderivación consiste en encontrar la última derivación a derechas que dió lugar a x . Esto es, si $x \in L(G)$ es porque existe una derivación a derechas de la forma

$$S \xRightarrow{*} yAz \Rightarrow ywz = x.$$

El problema es averiguar que regla $A \rightarrow w$ se aplicó y en que lugar de la cadena x se aplicó. En general, si queremos antiderivar una forma sentencial derecha $\beta\alpha w$ debemos averiguar por que regla $A \rightarrow \alpha$ seguir y en que lugar de la forma (después de β en el ejemplo) aplicarla.

$$S \xRightarrow{*} \beta Aw \Rightarrow \beta\alpha w.$$

La pareja formada por la regla y la posición se denomina mango o manecilla de la forma. Esta denominación viene de la visualización gráfica de la regla de producción como una mano que nos permite escalar hacia arriba en el árbol. Los “dedos” serían los símbolos en la parte derecha de la regla de producción.

Definición 7.2.2. Dada una gramática $G = (\Sigma, V, P, S)$ no ambigua, y dada una forma sentencial derecha $\alpha = \beta\gamma x$, con $x \in \Sigma^*$, el mango o handle de α es la última producción/posición que dió lugar a α :

$$\begin{array}{c} * \\ S \Longrightarrow \beta B x \Longrightarrow \beta \gamma x = \alpha \\ RM \end{array}$$

Escribiremos: $handle(\alpha) = (B \rightarrow \gamma, \beta\gamma)$. La función $handle$ tiene dos componentes: $handle_1(\alpha) = B \rightarrow \gamma$ y $handle_2(\alpha) = \beta\gamma$

Si dispusiéramos de un procedimiento que fuera capaz de identificar el mango, esto es, de detectar la regla y el lugar en el que se posiciona, tendríamos un mecanismo para construir un analizador. Lo curioso es que, a menudo es posible encontrar un autómata finito que reconoce el lenguaje de los prefijos $\beta\gamma$ que terminan en el mango. Con mas precisión, del lenguaje:

Definición 7.2.3. El conjunto de prefijos viables de una gramática G se define como el conjunto:

$$PV = \left\{ \delta \in (\Sigma \cup V)^* : \exists S \xRightarrow[RM]{*} \alpha \text{ y } \delta \text{ es un prefijo de } handle_2(\alpha) \right\}$$

Esto es, es el lenguaje de los prefijos viables es el conjunto de frases que son prefijos de $handle_2(\alpha)$ = $\beta\gamma$, siendo α una forma sentencial derecha ($\alpha \in FSD$). Los elementos de PV se denominan prefijos viables.

Obsérvese que si se dispone de un autómata que reconoce PV entonces se dispone de un mecanismo para investigar el lugar y el aspecto que pueda tener el mango. Si damos como entrada la sentencia $\alpha = \beta\gamma x$ a dicho autómata, el autómata aceptará la cadena $\beta\gamma$ pero rechazará cualquier extensión del prefijo. Ahora sabemos que el mango será alguna regla de producción de G cuya parte derecha sea un sufijo de $\beta\gamma$.

Definición 7.2.4. El siguiente autómata finito no determinista puede ser utilizado para reconocer el lenguaje de los prefijos viables PV :

- Alfabeto = $V \cup \Sigma$
- Los estados del autómata se denominan $LR(0)$ items. Son parejas formadas por una regla de producción de la gramática y una posición en la parte derecha de la regla de producción. Por ejemplo, $(E \rightarrow E + E, 2)$ sería un $LR(0)$ item para la gramática de las expresiones.

Conjunto de Estados:

$$Q = \{(A \rightarrow \alpha, n) : A \rightarrow \alpha \in P, n \leq |\alpha|\}$$

La notación $|\alpha|$ denota la longitud de la cadena $|\alpha|$. En vez de la notación $(A \rightarrow \alpha, n)$ escribiremos: $A \rightarrow \beta \uparrow \gamma = \alpha$, donde la flecha ocupa el lugar indicado por el número $n = |\beta|$:

- Función de transición:
- $$\delta(A \rightarrow \alpha \uparrow X \beta, X) = A \rightarrow \alpha X \uparrow \beta \quad \forall X \in V \cup \Sigma$$
- $$\delta(A \rightarrow \alpha \uparrow B \beta, \epsilon) = B \rightarrow \gamma \uparrow B \quad \forall B \in V$$
- Estado de arranque: Se añade la “superregla” $S' \rightarrow S$ a la gramática $G = (\Sigma, V, P, S)$. El $LR(0)$ item $S' \rightarrow \uparrow S$ es el estado de arranque.
 - Todos los estados definidos (salvo el de muerte) son de aceptación.

Denotaremos por $LR(0)$ a este autómata. Sus estados se denominan $LR(0) - items$. La idea es que este autómata nos ayuda a reconocer los prefijos viables PV .

Una vez que se tiene un autómata que reconoce los prefijos viables es posible construir un analizador sintáctico que construye una antiderivación a derechas. La estrategia consiste en “alimentar” el autómata con la forma sentencial derecha. El lugar en el que el autómata se detiene, rechazando indica el lugar exacto en el que termina el *handle* de dicha forma.

Ejemplo 7.2.1. Consideremos la gramática:

$$\begin{aligned} S &\rightarrow a S b \\ S &\rightarrow \epsilon \end{aligned}$$

El lenguaje generado por esta gramática es $L(G) = \{a^n b^n : n \geq 0\}$. Es bien sabido que el lenguaje $L(G)$ no es regular. La figura 7.1 muestra el autómata finito no determinista con ϵ -transiciones (NFA) que reconoce los prefijos viables de esta gramática, construido de acuerdo con el algoritmo 7.2.4.

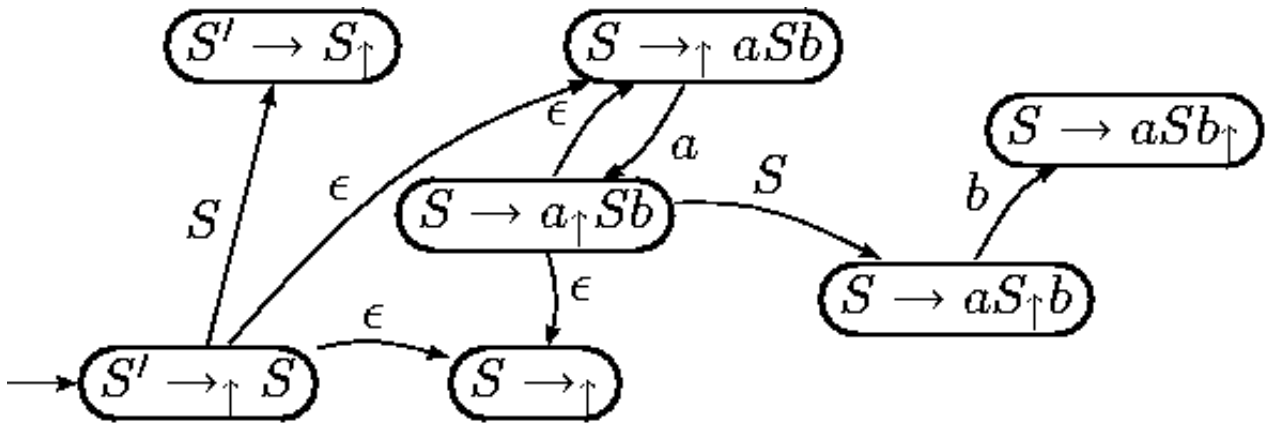


Figura 7.1: NFA que reconoce los prefijos viables

Ejercicio 7.2.1. Simule el comportamiento del autómata sobre la entrada $aabb$. ¿Donde rechaza? ¿En que estados está el autómata en el momento del rechazo?. ¿Qué etiquetas tienen? Haga también las trazas del autómata para las entradas $aaSbb$ y aSb . ¿Que antiderivación ha construido el autómata con sus sucesivos rechazos? ¿Que terminales se puede esperar que hayan en la entrada cuando se produce el rechazo del autómata?

7.3. Construcción de las Tablas para el Análisis SLR

7.3.1. Los conjuntos de Primeros y Siguietes

Repasemos las nociones de conjuntos de *Primeros* y *siguietes*:

Definición 7.3.1. Dada una gramática $G = (\Sigma, V, P, S)$ y un símbolo $\alpha \in (V \cup \Sigma)^*$ se define el conjunto $FIRST(\alpha)$ como:

$$FIRST(\alpha) = \{b \in \Sigma : \alpha \xRightarrow{*} b\beta\} \cup N(\alpha)$$

donde:

$$N(\alpha) = \begin{cases} \{\epsilon\} & \text{si } \alpha \xRightarrow{*} \epsilon \\ \emptyset & \text{en otro caso} \end{cases}$$

Definición 7.3.2. Dada una gramática $G = (\Sigma, V, P, S)$ y una variable $A \in V$ se define el conjunto $FOLLOW(A)$ como:

$$FOLLOW(A) = \left\{ b \in \Sigma : \exists S \xRightarrow{*} \alpha A b \right\} \cup E(A)$$

donde

$$E(A) = \begin{cases} \{\$ \} & \text{si } S \xRightarrow{*} \alpha A \\ \emptyset & \text{en otro caso} \end{cases}$$

Algoritmo 7.3.1. Construcción de los conjuntos $FIRST(X)$

1. Si $X \in \Sigma$ entonces $FIRST(X) = X$
2. Si $X \rightarrow \epsilon$ entonces $FIRST(X) = FIRST(X) \cup \{\epsilon\}$
3. Si $X \in V$ y $X \rightarrow Y_1 Y_2 \cdots Y_k \in P$ entonces

$i = 1;$

do

$FIRST(X) = FIRST(X) \cup FIRST(Y_i) - \{\epsilon\};$

$i++;$

mientras $(\epsilon \in FIRST(Y_i) \text{ and } (i \leq k))$

si $(\epsilon \in FIRST(Y_k) \text{ and } i > k)$ $FIRST(X) = FIRST(X) \cup \{\epsilon\}$

Este algoritmo puede ser extendido para calcular $FIRST(\alpha)$ para $\alpha = X_1 X_2 \cdots X_n \in (V \cup \Sigma)^*$.

Algoritmo 7.3.2. Construcción del conjunto $FIRST(\alpha)$

$i = 1;$

$FIRST(\alpha) = \emptyset;$

do

$FIRST(\alpha) = FIRST(\alpha) \cup FIRST(X_i) - \{\epsilon\};$

$i++;$

mientras $(\epsilon \in FIRST(X_i) \text{ and } (i \leq n))$

si $(\epsilon \in FIRST(X_n) \text{ and } i > n)$ $FIRST(\alpha) = FIRST(X) \cup \{\epsilon\}$

Algoritmo 7.3.3. Construcción de los conjuntos $FOLLOW(A)$ para las variables sintácticas $A \in V$:
Repetir los siguientes pasos hasta que ninguno de los conjuntos $FOLLOW$ cambie:

1. $FOLLOW(S) = \{\$ \}$ ($\$$ representa el final de la entrada)
2. Si $A \rightarrow \alpha B \beta$ entonces

$$FOLLOW(B) = FOLLOW(B) \cup (FIRST(\beta) - \{\epsilon\})$$

3. Si $A \rightarrow \alpha B$ o bien $A \rightarrow \alpha B \beta$ y $\epsilon \in FIRST(\beta)$ entonces

$$FOLLOW(B) = FOLLOW(B) \cup FOLLOW(A)$$

7.3.2. Construcción de las Tablas

Para la construcción de las tablas de un analizador SLR se construye el *autómata finito determinista* (DFA) (Q, Σ, δ, q_0) equivalente al NFA presentado en la sección 7.2 usando el *algoritmo de construcción del subconjunto*.

Como recordará, en la construcción del subconjunto, partiendo del estado de arranque q_0 del NFA con ϵ -transiciones se calcula su *clausura* $\overline{\{q_0\}}$ y las clausuras de los conjuntos de estados $\delta(\overline{\{q_0\}}, a)$ a

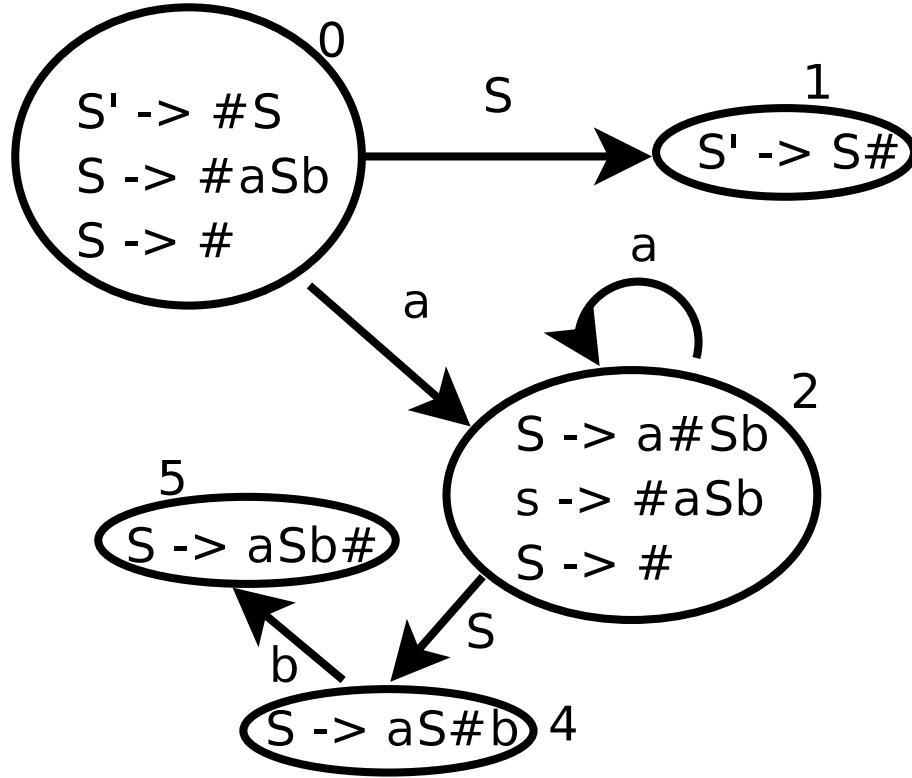


Figura 7.2: DFA equivalente al NFA de la figura 7.1

los que transita. Se repite el proceso con los conjuntos resultantes hasta que no se introducen nuevos conjuntos-estado.

La clausura \overline{A} de un subconjunto de estados del autómata A esta formada por todos los estados que pueden ser alcanzados mediante transiciones etiquetadas con la palabra vacía (denominadas ϵ transiciones) desde los estados de A . Se incluyen en \overline{A} , naturalmente los estados de A .

$$\overline{A} = \{q \in Q / \exists q' \in Q : \hat{\delta}(q, \epsilon) = q'\}$$

Aquí $\hat{\delta}$ denota la *función de transición del autómata* extendida a cadenas de Σ^* .

$$\hat{\delta}(q, x) = \begin{cases} \delta(\hat{\delta}(q, y), a) & \text{si } x = ya \\ q & \text{si } x = \epsilon \end{cases} \quad (7.1)$$

En la práctica, y a partir de ahora así lo haremos, se prescinde de diferenciar entre δ y $\hat{\delta}$ usándose indistintamente la notación δ para ambas funciones.

La clausura puede ser computada usando una estructura de pila o aplicando la expresión recursiva dada en la ecuación 7.1.

Para el NFA mostrado en el ejemplo 7.2.1 el DFA construido mediante esta técnica es el que se muestra en la figura 7.2. Se ha utilizado el símbolo $\#$ como marcador. Se ha omitido el número 3 para que los estados coincidan en numeración con los generados por **yapp** (véase el cuadro 7.1).

Un analizador sintáctico LR utiliza una tabla para su análisis. Esa tabla se construye a partir de la tabla de transiciones del DFA. De hecho, la tabla se divide en dos tablas, una llamada *tabla de saltos* o *tabla de gotos* y la otra *tabla de acciones*.

La tabla *goto* de un analizador *SLR* no es más que la tabla de transiciones del autómata DFA obtenido aplicando la construcción del subconjunto al NFA definido en 7.2.4. De hecho es la tabla de transiciones restringida a V (recuerde que el alfabeto del autómata es $V \cup \Sigma$). Esto es,

$$\delta|_{V \times Q} : V \times Q \rightarrow Q.$$

donde se define $goto(i, A) = \delta(A, I_i)$

La parte de la función de transiciones del DFA que corresponde a los terminales que no producen rechazo, esto es, $\delta|_{\Sigma \times Q} : \Sigma \times Q \rightarrow Q$ se adjunta a una tabla que se denomina *tabla de acciones*. La tabla de acciones es una tabla de doble entrada en los estados y en los símbolos de Σ . Las acciones de transición ante terminales se denominan *acciones de desplazamiento* o (*acciones shift*):

$$\delta|_{\Sigma \times Q} : \Sigma \times Q \rightarrow Q$$

donde se define $action(i, a) = \delta(a, I_i)$

Cuando un estado s contiene un LR(0)-item de la forma $A \rightarrow \alpha \uparrow$, esto es, el estado corresponde a un posible rechazo, ello indica que hemos llegado a un final del prefijo viable, que hemos visto α y que, por tanto, es probable que $A \rightarrow \alpha$ sea el *handle* de la forma sentencial derecha actual. Por tanto, añadiremos en entradas de la forma (s, a) de la tabla de acciones una acción que indique que hemos encontrado el mango en la posición actual y que la regla asociada es $A \rightarrow \alpha$. A una acción de este tipo se la denomina *acción de reducción*.

La cuestión es, ¿para que valores de $a \in \Sigma$ debemos disponer que la acción para (s, a) es de reducción? Podríamos decidir que ante cualquier terminal $a \in \Sigma$ que produzca un rechazo del autómata, pero podemos ser un poco mas selectivos. No cualquier terminal puede estar en la entrada en el momento en el que se produce la antiderivación o reducción. Observemos que si $A \rightarrow \alpha$ es el *handle* de γ es porque:

$$\begin{array}{ccc} * & & * \\ \exists S \Longrightarrow \beta A b x & \Longrightarrow & \beta \alpha b x = \gamma \\ RM & & RM \end{array}$$

Por tanto, cuando estamos reduciendo por $A \rightarrow \alpha$ los únicos terminales legales que cabe esperar en una reducción por $A \rightarrow \alpha$ son los terminales $b \in FOLLOW(A)$.

Dada una gramática $G = (\Sigma, V, P, S)$, podemos construir las tablas de acciones (*action table*) y transiciones (*gotos table*) mediante el siguiente algoritmo:

Algoritmo 7.3.4. *Construcción de Tablas SLR*

1. Utilizando el Algoritmo de Construcción del Subconjunto, se construye el Autómata Finito Determinista (DFA) $(Q, V \cup \Sigma, \delta, I_0, F)$ equivalente al Autómata Finito No Determinista (NFA) definido en 7.2.4. Sea $C = \{I_1, I_2, \dots, I_n\}$ el conjunto de estados del DFA. Cada estado I_i es un conjunto de LR(0)-items o estados del NFA. Asociemos un índice i con cada conjunto I_i .
2. La tabla de gotos no es más que la función de transición del autómata restringida a las variables de la gramática:

$$goto(i, A) = \delta(I_i, A) \text{ para todo } A \in V$$

3. Las acciones para el estado I_i se determinan como sigue:

- a) Si $A \rightarrow \alpha \uparrow a \beta \in I_i$, $\delta(I_i, a) = I_j$, $a \in \Sigma$ entonces:

$$action[i][a] = shift\ j$$

- b) Si $S' \rightarrow S \uparrow \in I_i$ entonces

$$action[i][\$] = accept$$

- c) Para cualquier otro caso de la forma $A \rightarrow \alpha \uparrow \in I_i$ distinto del anterior hacer

$$\forall a \in FOLLOW(A) : action[i][a] = reduce\ A \rightarrow \alpha$$

4. Las entradas de la tabla de acción que queden indefinidas después de aplicado el proceso anterior corresponden a acciones de “error”.

Definición 7.3.3. Si alguna de las entradas de la tabla resulta multievaluada, decimos que existe un conflicto y que la gramática no es SLR.

1. En tal caso, si una de las acciones es de ‘reducción’ y la otra es de ‘desplazamiento’, decimos que hay un conflicto shift-reduce o conflicto de desplazamiento-reducción.
2. Si las dos reglas indican una acción de reducción, decimos que tenemos un conflicto reduce-reduce o de reducción-reducción.

Ejemplo 7.3.1. Al aplicar el algoritmo 7.3.4 a la gramática 7.2.1

1	$S \rightarrow a S b$
2	$S \rightarrow \epsilon$

partiendo del autómata finito determinista que se construyó en la figura 7.2 y calculando los conjuntos de primeros y siguientes

	FIRST	FOLLOW
S	a, ϵ	$b, \$$

obtenemos la siguiente tabla de acciones SLR:

	a	b	$\$$
0	$s2$	$r2$	$r2$
1			aceptar
2	$s2$	$r2$	$r2$
4		$s5$	
5		$r1$	$r1$

Las entradas denotadas con $s\ n$ (s por shift) indican un desplazamiento al estado n , las denotadas con $r\ n$ (r por reduce o reducción) indican una operación de reducción o antiderivación por la regla n . Las entradas vacías corresponden a acciones de error.

El método de análisis *LALR* usado por **yapp** es una extensión del método SLR esbozado aquí. Supone un compromiso entre potencia (conjunto de gramáticas englobadas) y eficiencia (cantidad de memoria utilizada, tiempo de proceso). Veamos como **yapp** aplica la construcción del subconjunto a la gramática del ejemplo 7.2.1. Para ello construimos el siguiente programa **yapp**:

```
$ cat -n aSb.yp
1  %%
2  S:  # empty
3      |  'a' S 'b'
4  ;
5  %%
.....
```

y compilamos, haciendo uso de la opción **-v** para que **yapp** produzca las tablas en el fichero **aSb.output**:

```
$ ls -l aSb.*
-rw-r--r-- 1 lhp lhp 738 2004-12-19 09:52 aSb.output
-rw-r--r-- 1 lhp lhp 1841 2004-12-19 09:52 aSb.pm
-rw-r--r-- 1 lhp lhp 677 2004-12-19 09:46 aSb.yp
```

El contenido del fichero `aSb.output` se muestra en la tabla 7.1. Los números de referencia a las producciones en las acciones de reducción vienen dados por:

```
0: $start -> S $end
1: S -> /* empty */
2: S -> 'a' S 'b'
```

Observe que el final de la entrada se denota por `$end` y el marcador en un LR-item por un punto. Fíjese en el estado 2: En ese estado están también los items

`S -> . 'a' S 'b'` y `S -> .`

sin embargo no se explicitan por que se entiende que su pertenencia es consecuencia directa de aplicar la operación de clausura. Los LR items cuyo marcador no está al principio se denominan *items núcleo*.

Estado 0	Estado 1	Estado 2
<code>\$start -> . S \$end</code> <code>'a'shift 2</code> <code>\$default</code> <code>reduce 1 (S)</code> <code>S go to state 1</code>	<code>\$start -> S . \$end</code> <code>\$end shift 3</code>	<code>S -> 'a' . S 'b'</code> <code>'a'shift 2</code> <code>\$default</code> <code>reduce 1 (S)</code> <code>S go to state 4</code>
Estado 3	Estado 4	Estado 5
<code>\$start -> S \$end .</code> <code>\$default accept</code>	<code>S -> 'a' S . 'b'</code> <code>'b'shift 5</code>	<code>S -> 'a' S 'b' .</code> <code>\$default</code> <code>reduce 2 (S)</code>

Cuadro 7.1: Tablas generadas por `yapp`. El estado 3 resulta de transitar con `$`

Puede encontrar el listado completo de las tablas en `aSb.output` en el apéndice que se encuentra en la página ??.

Ejercicio 7.3.1. Compare la tabla 7.1 resultante de aplicar `yapp` con la que obtuvo en el ejemplo 7.3.1.

7.4. El módulo Generado por `yapp`

La ejecución de la orden `yapp -m Calc Calc.y` produce como salida el módulo `Calc.pm` el cual contiene las tablas LALR(1) para la gramática descrita en `Calc.y`. Estas tablas son las que dirigen al analizador LR. Puede ver el código completo del módulo en el apéndice que se encuentra en la página ??.

La estructura del módulo `Calc.pm` es como sigue:

```
1 package Calc;
2 use vars qw ( @ISA );
3 use strict;
4 @ISA= qw ( Parse::Yapp::Driver );
5 use Parse::Yapp::Driver;
6
```



```

7 sub new {
8     my($class)=shift;
9     ref($class) and $class=ref($class);
10
11     my($self)=$class->SUPER::new(
12         yyversion => '1.05',
13         yystates => [
14             ...
15             ], # estados
16         yyrules => [
17             # ... mas reglas
18             ], # final de las reglas
19         @_); # argumentos pasados
20     bless($self,$class);
21 }

```

La clase `Calc` hereda de `Parse::Yapp::Driver`, pero el objeto creado será bendecido en la clase `Calc` (Línea 4, véanse también la figura 7.3 y la línea 72 del fuente). Por tanto, el constructor llamado en la línea 11 es el de `Parse::Yapp::Driver`. Se utiliza la estrategia de llamada con parámetros con nombre. El valor para la clave `yystates` es una referencia anónima al array de estados y el valor para la clave `yyrules` es una referencia anónima a las reglas.

```

10
11     my($self)=$class->SUPER::new(
12         yyversion => '1.05',
13         yystates => [
14             {#State 0
15                 DEFAULT => -1, GOTOS => { 'input' => 1 }
16             },
17             {#State 1
18                 ACTIONS => {
19                     'NUM' => 6, ' ' => 4, "-" => 2, "(" => 7,
20                     'VAR' => 8, "\n" => 5, 'error' => 9
21                 },
22                 GOTOS => { 'exp' => 3, 'line' => 10 }
23             },
24             # ... mas estados
25             {#State 27
26                 ACTIONS => {
27                     "-" => 12, "+" => 13, "/" => 15, "^" => 16,
28                     "*" => 17
29                 },
30                 DEFAULT => -8
31             }
32         ], # estados

```

Se ve que un estado se pasa como un hash anónimo indexado en las acciones y los saltos.

Para consultar los números asociados con las reglas de producción vea el apéndice en la página ?? conteniendo el fichero `Calc.output`.

A continuación vienen las reglas:

```

33     yyrules => [
34         [#Rule 0
35             '$start', 2, undef ],
36         [#Rule 1

```

```

37         'input', 0, undef ],
38     [#Rule 2
39         'input', 2, sub
40 #line 17 "Calc.yp"
41         { push(@{$_[1]},$_[2]); $_[1] }
42     ],
43     [#Rule 3
44         'line', 1, sub
45 #line 20 "Calc.yp"
46         { $_[1] }
47     ],
48     [#Rule 4
49         'line', 2, sub
50 #line 21 "Calc.yp"
51         { print "$_[1]\n" }
52     ],
53 # ... mas reglas
54 [#Rule 11
55     'exp', 3, sub
56 #line 30 "Calc.yp"
57 { $_[1] * $_[3] }
58 ],
59 [#Rule 12
60     'exp', 3, sub
61 #line 31 "Calc.yp"
62 {
63     $_[3] and return($_[1] / $_[3]);
64     $_[0]->YYData->{ERRMSG} = "Illegal division by zero.\n";
65     $_[0]->YYError;
66     undef
67 }
68 ],
69 # ... mas reglas
70 ], # final de las reglas

```

Las reglas son arrays anónimos conteniendo el nombre de la regla o variable sintáctica (**exp**), el número de símbolos en la parte derecha y la subrutina anónima con el código asociado.

Vemos como la acción es convertida en una subrutina anónima. Los argumentos de dicha subrutina son los atributos semánticos asociados con los símbolos en la parte derecha de la regla de producción. El valor retornado por la acción/subrutina es el valor asociado con la reducción.

Para hacer que el compilador Perl diagnostique los errores relativos al fuente **Calc.yp** se usa una directiva **#line**.

```

71 @_;
72 bless($self,$class);
73 }
74

```

la bendición con dos argumentos hace que el objeto pertenezca a la clase **Calc**. A continuación siguen las subrutinas de soporte:

```

75 #line 44 "Calc.yp"
76
77
78 sub _Error {

```

```

79  # ...
80 }
81
82 sub _Lexer {
83   my($parser)=shift;
84   # ...
85 }
86
87 sub Run {
88   my($self)=shift;
89   $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
90 }
91
92 my($calc)=new Calc;
93 $calc->Run;
94
95 1;

```

7.5. Algoritmo de Análisis LR

Así pues la tabla de transiciones del autómata nos genera dos tablas: la tabla de acciones y la de saltos. El algoritmo de análisis sintáctico *LR* en el que se basa *yapp* utiliza una pila y dos tablas para analizar la entrada. Como se ha visto, la tabla de acciones contiene cuatro tipo de acciones:

1. Desplazar (*shift*)
2. Reducir (*reduce*)
3. Aceptar
4. Error

El algoritmo utiliza una pila en la que se guardan los estados del autómata. De este modo se evita tener que “comenzar” el procesado de la forma sentencial derecha resultante después de una reducción (antiderivación).

Algoritmo 7.5.1. *Análizador LR*

```

push(s0);
b = yylex();
for( ; ; ) {
  s = top(0); a = b;
  switch (action[s][a]) {
    case "shift t" :
      push(t);
      b = yylex();
      break;
    case "reduce A ->alpha" :
      eval(Sub{A -> alpha}->(top(|alpha|-1).attr, ... , top(0).attr));
      pop(|alpha|);
      push(goto[top(0)][A]);
      break;
    case "accept" : return (1);
    default : yyerror("syntax error");
  }
}
}

```

Como es habitual, $|x|$ denota la longitud de la cadena x . La función `top(k)` devuelve el elemento que ocupa la posición k desde el *top* de la pila (esto es, está a profundidad k). La función `pop(k)` extrae k elementos de la pila. La notación `state.attr` hace referencia al atributo asociado con cada estado. Denotamos por `sub_{reduce A -> alpha}` el código de la acción asociada con la regla $A \rightarrow \alpha$.

Todos los analizadores LR comparten, salvo pequeñas excepciones, el mismo algoritmo de análisis. Lo que más los diferencia es la forma en la que construyen las tablas. En `yapp` la construcción de las tablas de *acciones* y *gotos* se realiza mediante el algoritmo *LALR*.

7.6. Depuración en `yapp`

Es posible añadir un parámetro en la llamada a `YYParse` con nombre `yydebug` y valor el nivel de depuración requerido. Ello nos permite observar la conducta del analizador. Los posibles valores de depuración son:

Bit	Información de Depuración
0x01	Lectura de los terminales
0x02	Información sobre los estados
0x04	Acciones (shifts, reduces, accept ...)
0x08	Volcado de la pila
0x10	Recuperación de errores

Veamos un ejemplo de salida para la gramática que se describe en la página ?? cuando se llama con:

```
$self->YYParse( yylex => \&_Lexer, yyerror => \&_Error, yydebug => 0x1F )
```

```
1  $ ./use_aSb.pl
2  -----
3  In state 0:
4  Stack:[0]
5  ab # el usuario ha escrito esto
6  Need token. Got >a<
7  Shift and go to state 2.
8  -----
9  In state 2:
10 Stack:[0,2]
11 Need token. Got >b<
12 Reduce using rule 1 (S,0): S -> epsilon
13 Back to state 2, then go to state 4.
14 -----
15 In state 4:
16 Stack:[0,2,4]
17 Shift and go to state 5.
18 -----
19 In state 5:
20 Stack:[0,2,4,5]
21 Don't need token.
22 Reduce using rule 2 (S,3): S -> a S b
23 Back to state 0, then go to state 1.
24 -----
25 In state 1:
26 Stack:[0,1]
27 Need token. Got ><
```

```

28 Shift and go to state 3.
29 -----
30 In state 3:
31 Stack: [0,1,3]
32 Don't need token.
33 Accept.

```

7.7. Precedencia y Asociatividad

Recordemos que si al construir la tabla LALR, alguna de las entradas de la tabla resulta multievaluada, decimos que existe un conflicto. Si una de las acciones es de ‘reducción’ y la otra es de ‘desplazamiento’, se dice que hay un *conflicto shift-reduce* o *conflicto de desplazamiento-reducción*. Si las dos reglas indican una acción de reducción, decimos que tenemos un *conflicto reduce-reduce* o de *reducción-reducción*. En caso de que no existan indicaciones específicas *yapp* resuelve los conflictos que aparecen en la construcción de la tabla utilizando las siguientes reglas:

1. Un conflicto *reduce-reduce* se resuelve eligiendo la producción que se listó primero en la especificación de la gramática.
2. Un conflicto *shift-reduce* se resuelve siempre en favor del *shift*

Las declaraciones de precedencia y asociatividad mediante las palabras reservadas `%left` , `%right` , `%nonassoc` se utilizan para modificar estos criterios por defecto. La declaración de `token s` mediante la palabra reservada `%token` no modifica la precedencia. Si lo hacen las declaraciones realizadas usando las palabras `left` , `right` y `nonassoc` .

1. Los *tokens* declarados en la misma línea tienen igual precedencia e igual asociatividad. La precedencia es mayor cuanto mas abajo su posición en el texto. Así, en el ejemplo de la calculadora en la sección 7.1, el *token* `*` tiene mayor precedencia que `+` pero la misma que `/`.
2. La precedencia de una regla $A \rightarrow \alpha$ se define como la del terminal mas a la derecha que aparece en α . En el ejemplo, la producción

$$\text{expr} : \text{expr} \text{'+'} \text{expr}$$

tiene la precedencia del *token* `+`.

3. Para decidir en un conflicto *shift-reduce* se comparan la precedencia de la regla con la del terminal que va a ser desplazado. Si la de la regla es mayor se reduce si la del *token* es mayor, se desplaza.
4. Si en un conflicto *shift-reduce* ambos la regla y el terminal que va a ser desplazado tiene la misma precedencia *yapp* considera la asociatividad, si es asociativa a izquierdas, reduce y si es asociativa a derechas desplaza. Si no es asociativa, genera un mensaje de error.
Obsérvese que, en esta situación, la asociatividad de la regla y la del *token* han de ser por fuerza, las mismas. Ello es así, porque en *yapp* los *tokens* con la misma precedencia se declaran en la misma línea y sólo se permite una declaración por línea.
5. *Por tanto es imposible declarar dos tokens con diferente asociatividad y la misma precedencia.*
6. Es posible modificar la precedencia “natural” de una regla, calificándola con un *token* específico. para ello se escribe a la derecha de la regla `prec token`, donde `token` es un *token* con la precedencia que deseamos. Vea el uso del *token dummy* en el siguiente ejercicio.

Para ilustrar las reglas anteriores usaremos el siguiente programa *yapp*:

```
$ cat -n Precedencia.yp
 1 %token NUMBER
 2 %left '@'
 3 %right '&' dummy
 4 %%
 5 list
 6     :
 7     | list '\n'
 8     | list e
 9     ;
10
11 e : NUMBER
12   | e '&' e
13   | e '@' e %prec dummy
14   ;
15
16 %%
```

El código del programa cliente es el siguiente:

```
$ cat -n useprecedencia.pl
cat -n useprecedencia.pl
 1 #!/usr/bin/perl -w
 2 use strict;
 3 use Precedencia;
 4
 5 sub Error {
 6     exists $_[0]->YYData->{ERRMSG}
 7     and do {
 8         print $_[0]->YYData->{ERRMSG};
 9         delete $_[0]->YYData->{ERRMSG};
10         return;
11     };
12     print "Syntax error.\n";
13 }
14
15 sub Lexer {
16     my($parser)=shift;
17
18     defined($parser->YYData->{INPUT})
19     or $parser->YYData->{INPUT} = <STDIN>
20     or return('','undef');
21
22     $parser->YYData->{INPUT}=~s/^[ \t]//;
23
24     for ($parser->YYData->{INPUT}) {
25         s/^[0-9]+(?:\.[0-9]+)?//
26         and return('NUMBER',$1);
27         s/^(.)//s
28         and return($1,$1);
29     }
30 }
31
32 my $debug_level = (@ARGV)? oct(shift @ARGV): 0x1F;
```

```

33 my $parser = Precedencia->new();
34 $parser->YYParse( yylex => \&Lexer, yyerror => \&Error, yydebug => $debug_level );

```

Observe la llamada al analizador en la línea 34. Hemos añadido el parámetro con nombre *yydebug* con argumento `yydebug => $debug_level` (véase la sección 8.3 para ver los posibles valores de depuración).

Compilamos a continuación el módulo usando la opción `-v` para producir información sobre los conflictos y las tablas de salto y de acciones:

```

yapp -v -m Precedencia Precedencia.y
$ ls -ltr |tail -2
-rw-r--r--  1 lhp lhp   1628 2004-12-07 13:21 Precedencia.pm
-rw-r--r--  1 lhp lhp   1785 2004-12-07 13:21 Precedencia.output

```

La opción `-v` genera el fichero `Precedencia.output` el cual contiene información detallada sobre el autómata:

```

$ cat -n Precedencia.output
 1 Conflicts:
 2 -----
 3 Conflict in state 8 between rule 6 and token '@' resolved as reduce.
 4 Conflict in state 8 between rule 6 and token '&' resolved as shift.
 5 Conflict in state 9 between rule 5 and token '@' resolved as reduce.
 6 Conflict in state 9 between rule 5 and token '&' resolved as shift.
 7
 8 Rules:
 9 -----
10 0:      $start -> list $end
11 1:      list -> /* empty */
12 2:      list -> list '\n'
13 3:      list -> list e
14 4:      e -> NUMBER
15 5:      e -> e '&' e
16 6:      e -> e '@' e
17 ...

```

¿Porqué se produce un conflicto en el estado 8 entre la regla 6 (`e -> e '@' e`) y el terminal '@'? Editando el fichero `Precedencia.output` podemos ver los contenidos del estado 8:

```

85 State 8:
86
87      e -> e . '&' e (Rule 5)
88      e -> e . '@' e (Rule 6)
89      e -> e '@' e . (Rule 6)
90
91      '&'      shift, and go to state 7
92
93      $default      reduce using rule 6 (e)

```

El ítem de la línea 88 indica que debemos desplazar, el de la línea 89 que debemos reducir por la regla 6. ¿Porqué `yapp` resuelve el conflicto optando por reducir? ¿Que prioridad tiene la regla 6? ¿Que asociatividad tiene la regla 6? La declaración en la línea 13 modifica la precedencia y asociatividad de la regla:

```

13      | e '@' e %prec dummy

```

de manera que la regla pasa a tener la precedencia y asociatividad de `dummy`. Recuerde que habíamos declarado `dummy` como asociativo a derechas:

```
2 %left '@'
3 %right '&' dummy
```

¿Que ocurre? Que `dummy` tiene mayor prioridad que `'@'` y por tanto la regla tiene mayor prioridad que el terminal: por tanto se reduce.

¿Que ocurre cuando el terminal en conflicto es `'&'`? En ese caso la regla y el terminal tienen la misma prioridad. Se hace uso de la asociatividad a derechas que indica que el conflicto debe resolverse desplazando.

Ejercicio 7.7.1. *Explique la forma en que yacc resuelve los conflictos que aparecen en el estado 9. Esta es la información sobre el estado 9:*

State 9:

```
e -> e . '&' e (Rule 5)
e -> e '&' e . (Rule 5)
e -> e . '@' e (Rule 6)
'&'shift, and go to state 7
$default reduce using rule 5 (e)
```

Veamos un ejemplo de ejecución:

```
$ ./useprecedencia.pl
```

```
-----
In state 0:
Stack:[0]
Don't need token.
Reduce using rule 1 (list,0): Back to state 0, then go to state 1.
```

Lo primero que ocurre es una reducción por la regla en la que `list` produce vacío. Si miramos el estado 0 del autómata vemos que contiene:

```
20 State 0:
21
22 $start -> . list $end (Rule 0)
23
24 $default reduce using rule 1 (list)
25
26 list go to state 1
```

A continuación se transita desde 0 con `list` y se consume el primer terminal:

```
-----
In state 1:
Stack:[0,1]
2@3@4
Need token. Got >NUMBER<
Shift and go to state 5.
-----
In state 5:
Stack:[0,1,5]
Don't need token.
Reduce using rule 4 (e,1): Back to state 1, then go to state 2.
-----
```


En el estado 5 se reduce por la regla `e -> NUMBER`. Esto hace que se retire el estado 5 de la pila y se transite desde el estado 1 viendo el símbolo `e`:

```
In state 2:
Stack:[0,1,2]
Need token. Got >@<
Shift and go to state 6.
-----
In state 6:
Stack:[0,1,2,6]
Need token. Got >NUMBER<
Shift and go to state 5.
-----
In state 5:
Stack:[0,1,2,6,5]
Don't need token.
Reduce using rule 4 (e,1): Back to state 6, then go to state 8.
-----
In state 8:
Stack:[0,1,2,6,8]
Need token. Got >@<
Reduce using rule 6 (e,3): Back to state 1, then go to state 2.
-----
...
Accept.
```

Obsérvese la resolución del conflicto en el estado 8

La presencia de conflictos, aunque no siempre, en muchos casos es debida a la introducción de ambigüedad en la gramática. Si el conflicto es de desplazamiento-reducción se puede resolver explicitando alguna regla que rompa la ambigüedad. Los conflictos de reducción-reducción suelen producirse por un diseño erróneo de la gramática. En tales casos, suele ser mas adecuado modificar la gramática.

7.8. Generación interactiva de analizadores Yapp

En el siguiente código, la subrutina `create_yapp_package` nos muestra como crear un analizador Yapp en tiempo de ejecución. Las dos líneas:

```
my $p = new Parse::Yapp(input => $grammar);
$p = $p->Output(classname => $name);
```

crean una cadena en `$p` conteniendo el código de la clase que implanta el analizador. Todo el truco está en hacer

```
eval $p;
```

para tener el paquete a mano:

```
$ cat left.pl
#!/usr/local/bin/perl5.8.0 -w
#use strict;
use Parse::Yapp;

sub lex{
    my($parser)=shift;

    return('',undef) unless $parser->YYData->{INPUT};
```

```

        for ($parser->YYData->{INPUT}) {
            s/^\s*//;
            s/^(.)//;
            my $ret = $1;
            return($ret, $ret);
        }
    }

sub yapp {
    my $grammar = shift
        or die "Must specify a grammar as first argument";
    my $name = shift
        or die "Must specify the name of the class as second argument";

    my $p = new Parse::Yapp(input => $grammar) or die "Bad grammar.";
    $p = $p->Output(classname => $name) or die "Can't generate parser.";

    eval $p;
    $@ and die "Error while compiling your parser: $@\n";
}

##### main #####
my $grammar = q {
%left '*'
%%
S:  A
;

A:  A '*' A { "($_[1] $_[2] $_[3])" }
    | B
;

B:  'a' | 'b' | 'c' | 'd'
;

%%
};

&yapp($grammar, "Example");
my $p = new Example(yylex => \&lex, yyerror => sub {});

print "Expresion: ";
$p->YYData->{INPUT} = <>;
$p->YYData->{INPUT} =~ s/\s*$//;

my $out=$p->YYParse;
print "out = $out\n";

Sigue un ejemplo de ejecución:

$ ./left.pl
Expresion: a*b*c*d
out = (((a * b) * c) * d)

```

7.9. Construcción del Árbol Sintáctico

El siguiente ejemplo usa yacc para construir el árbol sintáctico de una expresión en infijo:

```
$ cat -n Infixtree_bless.y
1  #
2  # Infixtree.y
3  #
4
5  %{
6  use Data::Dumper;
7  %}
8  %right  '='
9  %left   '-' '+'
10 %left   '*' '/'
11 %left   NEG
12
13 %%
14 input:  #empty
15         |  input line
16 ;
17
18 line:    '\n'          { $_[1] }
19         | exp '\n'      { print Dumper($_[1]); }
20         | error '\n'    { $_[0]->YYError }
21 ;
22
23 exp:     NUM
24         |  VAR          { $_[1] }
25         |  VAR '=' exp   { bless $_[1], $_[3], 'ASSIGN' }
26         |  exp '+' exp   { bless $_[1], $_[3], 'PLUS' }
27         |  exp '-' exp   { bless $_[1], $_[3], 'MINUS' }
28         |  exp '*' exp   { bless $_[1], $_[3], 'TIMES' }
29         |  exp '/' exp   { bless $_[1], $_[3], 'DIVIDE' }
30         |  '-' exp %prec NEG { bless $_[2], 'NEG' }
31         |  '(' exp ')'    { $_[2] }
32 ;
33
34 %%
35
36 sub _Error {
37     exists $_[0]->YYData->{ERRMSG}
38     and do {
39         print $_[0]->YYData->{ERRMSG};
40         delete $_[0]->YYData->{ERRMSG};
41         return;
42     };
43     print "Syntax error.\n";
44 }
45
46 sub _Lexer {
47     my($parser)=shift;
48
49     defined($parser->YYData->{INPUT})
```

```

50     or $parser->YYData->{INPUT} = <STDIN>
51     or return('',undef);
52
53     $parser->YYData->{INPUT} =~ s/^[ \t]//;
54
55     for ($parser->YYData->{INPUT}) {
56         s/^[0-9]+(?:\.[0-9]+)?//
57         and return('NUM',$1);
58         s/^[A-Za-z][A-Za-z0-9_]*/
59         and return('VAR',$1);
60         s/^(.)//s
61         and return($1,$1);
62     }
63 }
64
65 sub Run {
66     my($self)=shift;
67     $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
68 }

```

Para compilar hacemos:

```
$ yacc -m Infixtree Infixtree_bless.y
```

El guión que usa el analizador anterior es similar al que vimos en la sección 7.1:

```

$ cat -n ./useinfixtree.pl
 1  #!/usr/bin/perl -w
 2
 3  use Infixtree;
 4
 5  $parser = new Infixtree();
 6  $parser->Run;

```

Veamos un ejemplo de ejecución:

```

$ ./useinfixtree.pl
a = 2+3
$VAR1 = bless( [
    'a',
    bless( [
        '2',
        '3',
    ], 'PLUS' ),
    ], 'ASSIGN' );
b = a*4+a
$VAR1 = bless( [
    'b',
    bless( [
        bless( [
            'a',
            '4',
        ], 'TIMES' ),
        'a',
    ], 'PLUS' ),
    ], 'ASSIGN' );

```

7.10. Acciones en Medio de una Regla

A veces necesitamos insertar una acción en medio de una regla. Una acción en medio de una regla puede hacer referencia a los atributos de los símbolos que la preceden (usando $\$n$), pero no a los que le siguen.

Cuando se inserta una acción $\{action_1\}$ para su ejecución en medio de una regla $A \rightarrow \alpha\beta$:

$$A \rightarrow \alpha \{action_1\} \beta \{action_2\}$$

yapp crea una variable sintáctica temporal T e introduce una nueva regla:

1. $A \rightarrow \alpha T \beta \{action_2\}$
2. $T \rightarrow \epsilon \{action_1\}$

Las acciones en mitad de una regla cuentan como un símbolo mas en la parte derecha de la regla. Así pues, en una acción posterior en la regla, se deberán referenciar los atributos de los símbolos, teniendo en cuenta este hecho.

Las acciones en mitad de la regla pueden tener un atributo. Las acciones posteriores en la regla se referirán a él como $\$_{[n]}$, siendo n su número de orden en la parte derecha.

7.11. Esquemas de Traducción

Un *esquema de traducción* es una gramática independiente del contexto en la cuál se han asociado atributos a los símbolos de la gramática. Un atributo queda caracterizado por un identificador o nombre y un tipo o clase. Además se han insertado acciones, esto es, código Perl/Python/C, ... en medio de las partes derechas. En ese código es posible referenciar los atributos de los símbolos de la gramática como variables del lenguaje subyacente.

Recuerde que el orden en que se evalúan los fragmentos de código es el de un recorrido primero-profundo del árbol de análisis sintáctico. Mas específicamente, considerando a las acciones como hijos-hoja del nodo, el recorrido que realiza un esquema de traducción es:

```
1  sub esquema_de_traducccion {
2      my $node = shift;
3
4      for my $child ($node->children) { # de izquierda a derecha
5          if ($child->isa('ACTION')) {
6              $child->execute;
7          }
8          else { esquema_de_traducccion($child) }
9      }
10 }
```

Obsérvese que, como el bucle de la línea 4 recorre a los hijos de izquierda a derecha, se debe dar la siguiente condición para que un esquema de traducción funcione:

Para cualquier regla de producción aumentada con acciones, de la forma

$$A \rightarrow X_1 \dots X_j \{ \text{action}(\$A\{b\}, \$X_1\{c\} \dots X_n\{d\}) \} X_{j+1} \dots X_n$$

debe ocurrir que los atributos evaluados en la acción insertada después de X_j dependan de atributos y variables que fueron computadas durante la visita de los hermanos izquierdos o de sus ancestros. En particular no deberían depender de atributos asociados con las variables $X_{j+1} \dots X_n$. Ello no significa que no sea correcto evaluar atributos de $X_{j+1} \dots X_n$ en esa acción.

7.12. Definición Dirigida por la Sintáxis

Una *definición dirigida por la sintáxis* es un pariente cercano de los esquemas de traducción. En una definición dirigida por la sintáxis una gramática $G = (V, \Sigma, P, S)$ se aumenta con nuevas características:

- A cada símbolo $S \in V \cup \Sigma$ de la gramática se le asocian cero o mas atributos. Un atributo queda caracterizado por un identificador o nombre y un tipo o clase. A este nivel son *atributos formales*, como los parámetros formales, en el sentido de que su realización se produce cuando el nodo del árbol es creado.
- A cada regla de producción $A \rightarrow X_1 X_2 \dots X_n \in P$ se le asocian un conjunto de *reglas de evaluación de los atributos* o *reglas semánticas* que indican que el atributo en la parte izquierda de la regla semántica depende de los atributos que aparecen en la parte derecha de la regla. El atributo que aparece en la parte izquierda de la regla semántica puede estar asociado con un símbolo en la parte derecha de la regla de producción.
- Los atributos de cada símbolo de la gramática $X \in V \cup \Sigma$ se dividen en dos grupos disjuntos: *atributos sintetizados* y *atributos heredados*. Un atributo de X es un *atributo heredado* si depende de atributos de su padre y hermanos en el árbol. Un *atributo sintetizado* es aquél tal que el valor del atributo depende de los valores de los atributos de los hijos, es decir en tal caso X ha de ser una variable sintáctica y los atributos en la parte derecha de la regla semántica deben ser atributos de símbolos en la parte derecha de la regla de producción asociada.
- Los atributos predefinidos se denominán *atributos intrínsecos*. Ejemplos de atributos intrínsecos son los atributos sintetizados de los terminales, los cuáles se han computado durante la fase de análisis léxico. También son atributos intrínsecos los atributos heredados del símbolo de arranque, los cuales son pasados como parámetros al comienzo de la computación.

La diferencia principal con un esquema de traducción está en que no se especifica el orden de ejecución de las reglas semánticas. Se asume que, bien de forma manual o automática, se resolverán las dependencias existentes entre los atributos determinadas por la aplicación de las reglas semánticas, de manera que serán evaluados primero aquellos atributos que no dependen de ningún otro, después los que dependen de estos, etc. siguiendo un esquema de ejecución que viene guiado por las dependencias existentes entre los datos.

Aunque hay muchas formas de realizar un evaluador de una definición dirigida por la sintáxis, conceptualmente, tal evaluador debe:

1. Construir el árbol de análisis sintáctico para la gramática y la entrada dadas.
2. Analizar las reglas semánticas para determinar los atributos, su clase y las dependencias entre los mismos.
3. Construir el *grafo de dependencias* de los atributos, el cual tiene un nodo para cada ocurrencia de un atributo en el árbol de análisis sintáctico etiquetado con dicho atributo. El grafo tiene una arista entre dos nodos si existe una dependencia entre los dos atributos a través de alguna regla semántica.
4. Supuesto que el grafo de dependencias determina un *orden parcial* (esto es cumple las propiedades reflexiva, antisimétrica y transitiva) construir un *orden topológico* compatible con el orden parcial.
5. Evaluar las reglas semánticas de acuerdo con el orden topológico.

Una definición dirigida por la sintáxis en la que las reglas semánticas no tienen efectos laterales se denomina una *gramática atribuída*.

Si la definición dirigida por la sintáxis puede ser realizada mediante un esquema de traducción se dice que es *L-atribuída*. Para que una definición dirigida por la sintáxis sea L-atribuída deben

cumplirse que cualquiera que sea la regla de producción $A \rightarrow X_1 \dots X_n$, los atributos heredados de X_j pueden depender únicamente de:

1. Los atributos de los símbolos a la izquierda de X_j
2. Los atributos heredados de A

Nótese que las restricciones se refieren a los atributos heredados. El cálculo de los atributos sintetizados no supone problema para un esquema de traducción. Si la gramática es LL(1), resulta fácil realizar una definición L-atribuida en un analizador descendente recursivo predictivo.

Si la definición dirigida por la sintaxis sólo utiliza atributos sintetizados se denomina *S-atribuida*. Una definición S-atribuida puede ser fácilmente trasladada a un programa **yapp**.

7.13. Manejo en yapp de Atributos Heredados

Supongamos que **yapp** esta inmerso en la construcción de la antiderivación a derechas y que la forma sentencial derecha en ese momento es:

$$X_m \dots X_1 X_0 Y_1 \dots Y_n a_1 \dots a_0$$

y que el mango es $B \rightarrow Y_1 \dots Y_n$ y en la entrada quedan por procesar $a_1 \dots a_0$.

Es posible acceder en **yapp** a los valores de los atributos de los estados en la pila del analizador que se encuentran “por debajo” o si se quiere “a la izquierda” de los estados asociados con la regla por la que se reduce. Para ello se usa una llamada al método **YYSemval**. La llamada es de la forma **\$_[0]->YYSemval(index)**, donde **index** es un entero. Cuando se usan los valores $1 \dots n$ devuelve lo mismo que **\$_[1]**, ..., **\$_[n]**. Esto es **\$_[1]** es el atributo asociado con Y_1 y **\$_[n]** es el atributo asociado con Y_n . Cuando se usa con el valor 0 devolverá el valor del atributo asociado con el símbolo que esta a la izquierda del mango actual, esto es el atributo asociado con X_0 , si se llama con -1 el que está dos unidades a la izquierda de la variable actual, esto es, el asociado con X_1 etc. Así **\$_[-m]** denota el atributo de X_m .

Esta forma de acceder a los atributos es especialmente útil cuando se trabaja con *atributos heredados*. Esto es, cuando un atributo de un nodo del árbol sintáctico se computa en términos de valores de atributos de su padre y/o sus hermanos. Ejemplos de atributos heredados son la clase y tipo en la declaración de variables. Supongamos que tenemos el siguiente *esquema de traducción* para calcular la clase (C) y tipo (T) en las declaraciones (D) de listas (L) de identificadores:

```

D → C T { $L{c} = $C{c}; $L{t} = $T{t} } L
C → global { $C{c} = "global" }
C → local { $C{c} = "local" }
T → integer { $T{t} = "integer" }
T → float { $T{t} = "float" }
L → { $L1{t} = $L{t}; $L1{c} = $L{c}; } L1 ','
    id { set_class($id{v}, $L{c}); set_type($id{v}, $L{t}); }
L → id { set_class($id{v}, $L{c}); set_type($id{v}, $L{t}); }

```

Los atributos **c** y **t** denotan respectivamente la clase y el tipo.

Ejercicio 7.13.1. *Evalúe el esquema de traducción para la entrada **global float x,y**. Represente el árbol de análisis, las acciones incrustadas y determine el orden de ejecución.*

*Olvide por un momento la notación usada en las acciones y suponga que se tratara de acciones **yapp**. ¿En que orden construye **yapp** el árbol y en que orden ejecutará las acciones?*

A la hora de transformar este esquema de traducción en un programa **yapp** es importante darse cuenta que en cualquier derivación a derechas desde D, cuando se reduce por una de las reglas

$$L \rightarrow id \mid L_1 \text{ '}' id$$

el símbolo a la izquierda de L es T y el que esta a la izquierda de T es C. Considere, por ejemplo la derivación a derechas:

$$\begin{aligned} D \Rightarrow C \ T \ L \Rightarrow C \ T \ L, id \Rightarrow C \ T \ L, id, id \Rightarrow C \ T \ id, id, id \Rightarrow \\ \Rightarrow C \ float \ id, id, id \Rightarrow local \ float \ id, id, id \end{aligned}$$

Observe que el orden de recorrido de **yapp** es:

$$\begin{aligned} local \ float \ id, id, id \Leftarrow C \ float \ id, id \Leftarrow C \ T \ id, id, id \Leftarrow \\ \Leftarrow C \ T \ L, id, id \Leftarrow C \ T \ L, id \Leftarrow C \ T \ L \Leftarrow D \end{aligned}$$

en la antiderivación, cuando el mango es una de las dos reglas para listas de identificadores, $L \rightarrow id$ y $L \rightarrow L, id$ es decir durante las tres ultimas antiderivaciones:

$$C \ T \ L, id, id \Leftarrow C \ T \ L, id \Leftarrow C \ T \ L \Leftarrow D$$

las variables a la izquierda del mango son T y C. Esto ocurre siempre. Estas observaciones nos conducen al siguiente programa **yapp**:

```
$ cat -n Inherited.y
1  %token FLOAT INTEGER
2  %token GLOBAL
3  %token LOCAL
4  %token NAME
5
6  %%
7  declarationlist
8      : # vacio
9      | declaration ';' declarationlist
10     ;
11
12  declaration
13      : class type namelist { ; }
14      ;
15
16  class
17      : GLOBAL
18      | LOCAL
19      ;
20
21  type
22      : FLOAT
23      | INTEGER
24      ;
25
26  namelist
27      : NAME
28      { printf("%s de clase %s, tipo %s\n",
29              $_[1], $_[0]->YYSemval(-1), $_[0]->YYSemval(0)); }
30      | namelist ',' NAME
31      { printf("%s de clase %s, tipo %s\n",
32              $_[3], $_[0]->YYSemval(-1), $_[0]->YYSemval(0)); }
33      ;
34  %%
```


A continuación escribimos el programa que usa el módulo generado por yapp:

```
$ cat -n useinherited.pl
1  #!/usr/bin/perl -w
2  use strict;
3  use Inherited;
4
5  sub Error {
6      exists $_[0]->YYData->{ERRMSG}
7      and do {
8          print $_[0]->YYData->{ERRMSG};
9          delete $_[0]->YYData->{ERRMSG};
10         return;
11     };
12     print "Error sintáctico\n";
13 }
14
15 { # hagamos una clausura con la entrada
16     my $input;
17     local $/ = undef;
18     print "Entrada (En Unix, presione CTRL-D para terminar):\n";
19     $input = <stdin>;
20
21     sub scanner {
22
23         { # Con el redo del final hacemos un bucle "infinito"
24             if ($input =~ m|\G\s*INTEGER\b|igc) {
25                 return ('INTEGER', 'INTEGER');
26             }
27             elsif ($input =~ m|\G\s*FLOAT\b|igc) {
28                 return ('FLOAT', 'FLOAT');
29             }
30             elsif ($input =~ m|\G\s*LOCAL\b|igc) {
31                 return ('LOCAL', 'LOCAL');
32             }
33             elsif ($input =~ m|\G\s*GLOBAL\b|igc) {
34                 return ('GLOBAL', 'GLOBAL');
35             }
36             elsif ($input =~ m|\G\s*([a-z_]\w*)\b|igc) {
37                 return ('NAME', $1);
38             }
39             elsif ($input =~ m|\G\s*([,;])|gc) {
40                 return ($1, $1);
41             }
42             elsif ($input =~ m|\G\s*(.)|gc) {
43                 die "Caracter invalido: $1\n";
44             }
45             else {
46                 return ('', undef); # end of file
47             }
48             redo;
49         }
50     }
51 }
```

```

52
53 my $debug_level = (@ARGV)? oct(shift @ARGV): 0x1F;
54 my $parser = Inherited->new();
55 $parser->YYParse( yylex => \&scanner, yyerror => \&Error, yydebug => $debug_level );

```

En las líneas de la 15 a la 51 está nuestro analizador léxico. La entrada se lee en una variable local cuyo valor permanece entre llamadas: hemos creado una clausura con la variable `$input` (véase la sección [?] para más detalles sobre el uso de clausuras en Perl). Aunque la variable `$input` queda inaccesible desde fuera de la clausura, persiste entre llamadas como consecuencia de que la subrutina `scanner` la utiliza.

A continuación sigue un ejemplo de ejecución:

```

$ ./useinherited.pl 0
Entrada (En Unix, presione CTRL-D para terminar):
global integer x, y, z;
local float a,b;
x de clase GLOBAL, tipo INTEGER
y de clase GLOBAL, tipo INTEGER
z de clase GLOBAL, tipo INTEGER
a de clase LOCAL, tipo FLOAT
b de clase LOCAL, tipo FLOAT

```

Ejercicio 7.13.2. *El siguiente programa yapp calcula un árbol de análisis abstracto para la gramática del ejemplo anterior:*

```

%token FLOAT INTEGER
%token GLOBAL
%token LOCAL
%token NAME

%%
declarationlist
: /* vacio */ { bless [], 'declarationlist' }
| declaration ';' declarationlist { push @{$_[3]}, $_[1]; $_[3] }
;

declaration
: class type namelist
{
    bless {class => $_[1], type => $_[2], namelist => $_[3]}, 'declaration';
}
;

class
: GLOBAL { bless { GLOBAL => 0}, 'class' }
| LOCAL { bless { LOCAL => 1}, 'class' }
;

type
: FLOAT { bless { FLOAT => 2}, 'type' }
| INTEGER { bless { INTEGER => 3}, 'type' }
;

namelist

```

```

: NAME
  { bless [ $_[1]], 'namelist' }
| namelist ', ' NAME
  { push @{$_[1]}, $_[3]; $_[1] }
;
%%

```

sigue un ejemplo de ejecución:

```

$ ./useinherited3.pl
Entrada (En Unix, presione CTRL-D para terminar):
global float x,y;
$VAR1 = bless( [
  bless( {
    'namelist' => bless( [ 'x', 'y' ], 'namelist' ),
    'type' => bless( { 'FLOAT' => 2 }, 'type' ),
    'class' => bless( { 'GLOBAL' => 0 }, 'class' )
  }, 'declaration' )
], 'declarationlist' );

```

Extienda el programa del ejemplo para que la gramática incluya las acciones del esquema de traducción. Las acciones se tratarán como un terminal `CODE` y serán devueltas por el analizador léxico. Su atributo asociado es el texto del código. El programa `yapp` deberá devolver el árbol abstracto extendido con las acciones-terminales. La parte mas difícil de este problema consiste en “reconocer” el código Perl incrustado. La estrategia seguir consiste en contar el número de llaves que se abren y se cierran. Cuando el contador alcanza cero es que hemos llegado al final del código Perl incrustado. Esta estrategia tiene una serie de problemas. ¿Sabría decir cuáles? (sugerencia: repase la sección 7.19.3 o vea como `yapp` resuelve el problema).

7.14. Acciones en Medio de una Regla y Atributos Heredados

La estrategia utilizada en la sección 8.28 funciona si podemos predecir la posición del atributo en la pila del analizador. En el ejemplo anterior los atributos clase y tipo estaban siempre, cualquiera que fuera la derivación a derechas, en las posiciones 0 y -1. Esto no siempre es así. Consideremos la siguiente *definición dirigida por la sintaxis*:

$S \rightarrow a A C$	$\$C\{i\} = \$A\{s\}$
$S \rightarrow b A B C$	$\$C\{i\} = \$A\{s\}$
$C \rightarrow c$	$\$C\{s\} = \$C\{i\}$
$A \rightarrow a$	$\$A\{s\} = "a"$
$B \rightarrow b$	$\$B\{s\} = "b"$

Ejercicio 7.14.1. *Determine un orden correcto de evaluación de la anterior definición dirigida por la sintaxis para la entrada `b a b c`.*

C hereda el atributo sintetizado de A. El problema es que, en la pila del analizador el atributo $\$A\{s\}$ puede estar en la posición 0 o -1 dependiendo de si la regla por la que se derivó fue $S \rightarrow a A C$ o bien $S \rightarrow b A B C$. La solución a este tipo de problemas consiste en insertar acciones intermedias de copia del atributo de manera que se garantice que el atributo de interés está siempre a una distancia fija. Esto es, se inserta una variable sintáctica intermedia auxiliar M la cual deriva a vacío y que tiene como acción asociada una regla de copia:

$S \rightarrow a A C$	$\$C\{i\} = \$A\{s\}$
$S \rightarrow b A B M C$	$\$M\{i\} = \$A\{s\}; \$C\{i\} = \$M\{s\}$
$C \rightarrow c$	$\$C\{s\} = \$C\{i\}$
$A \rightarrow a$	$\$A\{s\} = "a"$
$B \rightarrow b$	$\$B\{s\} = "b"$
$M \rightarrow \epsilon$	$\$M\{s\} = \$M\{i\}$

El nuevo esquema de traducción puede ser implantado mediante un programa `yapp`:

```
$ cat -n Inherited2.y
1  %%
2  S : 'a' A C
3    | 'b' A B { $_[2]; } C
4    ;
5
6  C : 'c' { print "Valor: ",$_[0]->YYSemval(0),"\\n"; $_[0]->YYSemval(0) }
7    ;
8
9  A : 'a' { 'a' }
10   ;
11
12 B : 'b' { 'b' }
13   ;
14
15 %%
```

La ejecución muestra como se ha propagado el valor del atributo:

```
$ ./useinherited2.pl '0x04'
Entrada (En Unix, presione CTRL-D para terminar):
b a b c
Shift 2. Shift 6.
Reduce using rule 5 (A,1): Back to state 2, then state 5.
Shift 8.
Reduce 6 (B,1): Back to state 5, then state 9.
Reduce 2 (@1-3,0): Back to state 9, then state 12.
```

En este momento se esta ejecutando la acción intermedia. Lo podemos comprobar revisando el fichero `Inherited2.output` que fué generado usando la opción `-v` al llamar a `yapp`. La regla 2 por la que se reduce es la asociada con la acción intermedia:

```
$ cat -n Inherited2.output
1 Rules:
2 -----
3 0:      $start -> S $end
4 1:      S -> 'a' A C
5 2:      @1-3 -> /* empty */
6 3:      S -> 'b' A B @1-3 C
7 4:      C -> 'c'
8 5:      A -> 'a'
9 6:      B -> 'b'
...

```

Obsérvese la notación usada por `yapp` para la *acción en medio de la regla*: `@1-3`. Continuamos con la antiderivación:

Shift 10.
 Reduce 4 (C,1):
 Valor: a
 Back to state 12, then 13.
 Reduce using rule 3 (S,5): Back to state 0, then state 1.
 Shift 4.
 Accept.

El método puede ser generalizado a casos en los que el atributo de interés este a diferentes distancias en diferentes reglas sin mas que introducir las correspondientes acciones intermedias de copia.

7.15. Recuperación de Errores

Las entradas de un traductor pueden contener errores. El lenguaje **yapp** proporciona un *token* especial, **error**, que puede ser utilizado en el programa fuente para extender el traductor con “producciones de error” que lo doten de cierta capacidad para recuperarse de una entrada errónea y poder continuar analizando el resto de la entrada.

Consideremos lo que ocurre al ejecutar nuestra calculadora **yapp** con una entrada errónea. Recordemos la gramática:

```

 9 %right  '='
10 %left  '- ' '+'
11 %left  '* ' '/'
12 %left  NEG
13 %right  '^'
14
15 %%
16 input:  # empty
17         |  input line  { push(@{$_[1]},$_[2]); $_[1] }
18 ;
19
20 line:    '\n'          { $_[1] }
21         |  exp '\n'    { print "$_[1]\n" }
22         |  error '\n'  { $_[0]->YYErrork }
23 ;

```

La regla `line → error '\n'` es una producción de error. La idea general de uso es que, a través de la misma, el programador le indica a **yapp** que, cuando se produce un error dentro de una expresión, descarte todos los *tokens* hasta llegar al retorno del carro y prosiga con el análisis. Además, mediante la llamada al método **YYErrork** el programador anuncia que, si se alcanza este punto, la recuperación puede considerarse “completa” y que **yapp** puede emitir a partir de ese momento mensajes de error con la seguridad de que no son consecuencia de un comportamiento inestable provocado por el primer error.

El resto de la gramática de la calculadora era como sigue:

```

24
25 exp:      NUM
26         |  VAR          { $_[0]->YYData->{VARS}{$_[1]} }
27         |  VAR '=' exp  { $_[0]->YYData->{VARS}{$_[1]}=$_[3] }
28         |  exp '+' exp  { $_[1] + $_[3] }
29         |  exp '-' exp  { $_[1] - $_[3] }
30         |  exp '*' exp  { $_[1] * $_[3] }
31         |  exp '/' exp  {
32                        $_[3]

```

```

33                                     and return($_[1] / $_[3]);
34                                     $_[0]->YYData->{ERRMSG}
35                                     = "Illegal division by zero.\n";
36                                     $_[0]->YYError;
37                                     undef
38                                     }
39         | '-' exp %prec NEG { -$_[2] }
40         | exp '^' exp { $_[1] ** $_[3] }
41         | '(' exp ')' { $_[2] }
42 ;

```

en la ejecución activamos el flag `yydebug` a 0x10 para obtener información sobre el tratamiento de errores:

```
$self->YYParse( yylex => \&_Lexer, yyerror => \&_Error, yydebug => 0x10 );
```

Pasemos a darle una primera entrada errónea:

```

$ ./usecalc.pl
3-+2
Syntax error.
**Entering Error recovery.
**Pop state 12.
**Pop state 3.
**Shift $error token and go to state 9.
**Dicard invalid token >+<.
**Pop state 9.
**Shift $error token and go to state 9.
**Dicard invalid token >NUM<.
**Pop state 9.
**Shift $error token and go to state 9.
**End of Error recovery.

```

El esquema general del algoritmo de recuperación de errores usado por la versión actual de `yapp` es el siguiente:

1. Cuando se encuentra ante una acción de error, el analizador genera un *token error*.
2. A continuación pasa a retirar estados de la pila hasta que descubre un estado capaz de transitar ante el *token error*. En el ejemplo anterior el analizador estaba en el estado 12 y lo retira de la pila. Los contenidos del estado 12 son:

```

exp -> exp '-' . exp (Rule 10)
'('shift 7 '-'shift 2 NUM shift 6 VAR shift 8
exp go to state 21

```

Obviamente no esperabamos ver un `'+'` aquí. El siguiente estado en la cima de la pila es el 3, el cual tampoco tiene ninguna transición ante el terminal **error**:

```

line -> exp . '\n'(Rule 4)
exp -> exp . '+' exp (Rule 9)
exp -> exp . '-' exp (Rule 10)
exp -> exp . '*' exp (Rule 11)
exp -> exp . '/' exp (Rule 12)
exp -> exp . '^' exp (Rule 14)

'*'shift 17 '+'shift 13 '-'shift 12 '/'shift 15
'\n'shift 14 '^'shift 16

```

El pop sobre el estado 3 deja expuesto en la superficie al estado 1, el cuál “sabe” como manejar el error:

```
$start -> input . $end (Rule 0)
input -> input . line (Rule 2)

$end shift 4 '('shift 7 '-'shift 2 '\n'shift 5
NUM shift 6 VAR shift 8
error shift 9

exp go to state 3
line go to state 10
```

3. En este punto transita al estado correspondiente a desplazar el *token error*.

En consecuencia, con lo dicho, en el ejemplo se va al estado 9:

```
line -> error . '\n'(Rule 5)
'\n'shift, and go to state 20
```

4. Entonces el algoritmo de recuperación va leyendo *tokens* y descartandolos hasta encontrar uno que sea aceptable. En este caso hemos especificado que el terminal que nos da cierta confianza de recuperación es el retorno de carro:

```
**Discard invalid token >+<.
**Pop state 9.
**Shift $error token and go to state 9.
**Discard invalid token >NUM<.
**Pop state 9.
**Shift $error token and go to state 9.
**End of Error recovery.
```

5. Sólo se envían nuevos mensajes de error una vez asimilados (desplazados) algunos símbolos terminales. De este modos se intenta evitar la aparición masiva de mensajes de error.

7.16. Recuperación de Errores en Listas

Aunque no existe un método exacto para decidir como ubicar las reglas de recuperación de errores, en general, los símbolos de error deben ubicarse intentado satisfacer las siguientes reglas:

- Tan cerca como sea posible del símbolo de arranque.
- Tan cerca como sea posible de los símbolos terminales.
- Sin introducir nuevos conflictos.

En el caso particular de las listas, se recomienda seguir el siguiente esquema:

Ejercicio 7.16.1. Compruebe el funcionamiento de la metodología para la recuperación de errores en listas presentada en la tabla 7.2 estudie el siguiente programa **yapp** siguiendo la traza de estados, generando entradas con todos los tipos de error posibles. ¿Cómo se recupera el analizador en caso de existencia de un segundo error? ¿Que ocurre si dos errores consecutivos están muy próximos? El programa corresponde al tercer caso de la tabla 7.2, el caso $x:y\{Ty\}$ con $x = \text{list}$, $T = ', '$ e $y = \text{NUMBER}$:

Construcción	EBNF	yapp
secuencia opcional	$x:\{y\}$	<pre>x : /* null */ x y { \$_[0]->YYError; } x error</pre>
secuencia	$x:y\{y\}$	<pre>x : y xy { \$_[0]->YYError; } error x error</pre>
lista	$x:y\{Ty\}$	<pre>x : y x T y { \$_[0]->YYError; } error x error x error y { \$_[0]->YYError; } x T error</pre>

Cuadro 7.2: Recuperación de errores en listas

```
%token NUMBER
%%
command
:
| command list '\n' { $_[0]->YYError; }
;

list
: NUMBER          { put($1); }
| list ',' NUMBER { put($3); $_[0]->YYError; }
| error           { err(1); }
| list error      { err(2); }
| list error NUMBER { err(3); put($3); $_[0]->YYError; }
| list ',' error  { err(4); }
;

%%
sub put { my $x = shift; printf("%2.1lf\n",$x); }
sub err { my $code = shift; printf("err %d\n",$code); }
...
```

7.17. Consejos a seguir al escribir un programa yapp

Cuando escriba un programa **yapp** asegurese de seguir los siguientes consejos:

1. Coloque el punto y coma de separación de reglas en una línea aparte. Un punto y coma “pegado” al final de una regla puede confundirse con un terminal de la regla.
2. Si hay una regla que produce vacío, coloquela en primer lugar y acompáñela de un comentario resaltando ese hecho.
3. Nunca escriba dos reglas de producción en la misma línea.
4. Sangre convenientemente todas las partes derechas de las reglas de producción de una variable, de modo que queden alineadas.

5. Ponga nombres representativos a sus variables sintácticas. No llame *Z* a una variable que representa el concepto “lista de parámetros”, llámela `ListaDeParametros`.
6. Es conveniente que declare los terminales simbólicos, esto es, aquellos que llevan un identificador asociado. Si no llevan prioridad asociada o no es necesaria, use una declaración `%token`. De esta manera el lector de su programa se dará cuenta rápidamente que dichos identificadores no se corresponden con variables sintácticas. Por la misma razón, si se trata de terminales asociados con caracteres o cadenas no es tan necesario que los declare, a menos que, como en el ejemplo de la calculadora para `'+'` y `'*'`, sea necesario asociarles una precedencia.
7. Es importante que use la opción `-v` para producir el fichero `.output` conteniendo información detallada sobre los conflictos y el autómata. Cuando haya un conflicto shift-reduce no resuelto busque en el fichero el estado implicado y vea que LR(0) items $A \rightarrow \alpha \uparrow$ y $B \rightarrow \beta \uparrow \gamma$ entran en conflicto.
8. Si según el informe de `yapp` el conflicto se produce ante un terminal a , es porque $a \in FOLLOW(A)$ y $a \in FIRST(\gamma)$. Busque las causas por las que esto ocurre y modifique su gramática con vistas a eliminar la presencia del terminal a en uno de los dos conjuntos implicados o bien establezca reglas de prioridad entre los terminales implicados que resuelvan el conflicto.
9. Nótese que cuando existe un conflicto de desplazamiento reducción entre $A \rightarrow \alpha \uparrow$ y $B \rightarrow \beta \uparrow \gamma$, el programa `yapp` contabiliza un error por cada terminal $a \in FOLLOW(A) \cap FIRST(\gamma)$. Por esta razón, si hay 16 elementos en $FOLLOW(A) \cap FIRST(\gamma)$, el analizador `yapp` informará de la existencia de 16 conflictos *shift-reduce*, cuando en realidad se trata de uno sólo. No desespere, los conflictos “auténticos” suelen ser menos de los que `yapp` anuncia.
10. Si necesita declarar variables globales, inicializaciones, etc. que afectan la conducta global del analizador, escriba el código correspondiente en la cabecera del analizador, protegido por los delimitadores `%{` y `%}`. Estos delimitadores deberán aparecer en una línea aparte. Por ejemplo:

```
%{
our contador = 0;
}%

%token NUM
...
%%
```

11. Si tiene problemas en tiempo de ejecución con el comportamiento del analizador sintáctico use la opción `yydebug => 0x1F` en la llamada al analizador.
12. Si trabaja en windows y pasa los ficheros a unix tenga cuidado con la posible introducción de caracteres espúreos en el fichero. Debido a la presencia de caracteres de control invisibles, el analizador `yapp` pasará a rechazar una gramática aparentemente correcta.
13. Sea consciente de que los analizadores sintáctico y léxico mantienen una relación de corutinas en `yapp`: Cada vez que el analizador sintáctico necesita un nuevo terminal para decidir que regla de producción se aplica, llama al analizador léxico, el cuál deberá retornar el siguiente terminal. La estrategia es diferente de la utilizada en el ejemplo usado para el lenguaje Tutu en el capítulo 4. Allí generábamos en una primera fase la lista de terminales. Aquí los terminales se generan de uno en uno y cada vez que se encuentra uno nuevo se retorna al analizador sintáctico. La ventaja que tiene este método es que permite colaborar al analizador sintáctico y al analizador léxico para “dinámicamente” modificar la conducta del análisis léxico. Por ejemplo en los compiladores del lenguaje C es común hacer que el analizador léxico cuando descubre un identificador que previamente ha sido declarado como identificador de tipo (mediante el uso de `typedef`) retorne un terminal `TYPENAME` diferente del terminal `ID` que caracteriza a los identificadores. Para ello,

el analizador sintáctico, cuando detecta una tal declaración, “avisa” al analizador léxico para que modifique su conducta. El analizador sintáctico volverá a avisarlo cuando la declaración del identificador como identificador de tipo salga de ámbito y pierda su especial condición.

14. En **yapp** el analizador sintáctico espera que el analizador léxico devuelva de cada vez una pareja formada por dos escalares. El primer escalar es la cadena que designa el terminal. A diferencia de la habitual costumbre **yacc** de codificar los terminales como enteros, en **yapp** se suelen codificar como cadenas. La segunda componente de la pareja es el atributo asociado con el terminal. Si el atributo es un atributo complejo que necesitas representar mediante un hash o un vector, lo mejor es hacer que esta componente sea una referencia al objeto describiendo el atributo. El analizador léxico le indica al sintáctico la finalización de la entrada enviándole la pareja (`''`, `undef`) formada por la palabra vacía con atributo `undef`.
15. Hay fundamentalmente dos formas de hacer el analizador léxico: hacerlo destructivo o no destructivo. En los destructivos se usa el operador de sustitución `s` (véase el ejemplo de la sección 7.1), en cuyo caso la entrada procesada es retirada de la cadena leída. En los no destructivos utilizamos el operador de emparejamiento `m`. Véase el ejemplo de analizador léxico en la sección 8.28 (concretamente la subrutina `scanner` en la línea 20 del fichero `useinherited.pl`)

Ejemplo 7.17.1. *Consideremos de nuevo el programa **yapp** para producir árboles para las expresiones en infijo. Supongamos que olvidamos introducir una prioridad explícita al terminal '=':*

```
$ cat -n Infixtree_conflict.yp
1  #
2  # Infixtree.yp
3  #
4
5  %{
6  use Data::Dumper;
7  %}
8  %left  '-' '+'
9  %left  '*' '/'
10 %left  NEG
11
12 %%
13 input:  #empty
14         | input line
15 ;
16
17 line:    '\n'          { $_[1] }
18         | exp '\n'      { print Dumper($_[1]); }
19         | error '\n'    { $_[0]->YYError }
20 ;
21
22 exp:     NUM
23         | VAR           { $_[1] }
24         | VAR '=' exp    { bless $_[1], $_[3], 'ASSIGN' }
25         | exp '+' exp    { bless $_[1], $_[3], 'PLUS' }
26         | exp '-' exp    { bless $_[1], $_[3], 'MINUS' }
27         | exp '*' exp    { bless $_[1], $_[3], 'TIMES' }
28         | exp '/' exp    { bless $_[1], $_[3], 'DIVIDE' }
29
30 ....
```

en este caso al compilar encontraremos conflictos:

```
$ yacc -v -m Infixtree Infixtree_conflict.y
4 shift/reduce conflicts
```

En tal caso lo que debemos hacer es editar el fichero .output. El comienzo del fichero es como sigue:

```
$ cat -n Infixtree_conflict.output
1 Warnings:
2 -----
3 4 shift/reduce conflicts
4
5 Conflicts:
6 -----
7 Conflict in state 11 between rule 13 and token '-' resolved as reduce.
8 Conflict in state 11 between rule 13 and token '*' resolved as reduce.
...

```

Tal y como indica la expresión ...resolved as ..., las líneas como la 7, la 8 y siguientes se refieren a conflictos resueltos. Mas abajo encontraremos información sobre la causa de nuestros conflictos no resueltos:

```
...
26 Conflict in state 23 between rule 11 and token '/' resolved as reduce.
27 State 25 contains 4 shift/reduce conflicts

```

Lo que nos informa que los conflictos ocurren en el estado 25 ante 4 terminales distintos. Nos vamos a la parte del fichero en la que aparece la información relativa al estado 25. Para ello, como el fichero es grande, buscamos por la cadena adecuada. En vi buscaríamos por /^State 25. Las líneas correspondientes contienen:

```
291 State 25:
292
293 exp -> VAR '=' exp . (Rule 8)
294 exp -> exp . '+' exp (Rule 9)
295 exp -> exp . '-' exp (Rule 10)
296 exp -> exp . '*' exp (Rule 11)
297 exp -> exp . '/' exp (Rule 12)
298
299 '*' shift, and go to state 16
300 '+' shift, and go to state 13
301 '-' shift, and go to state 12
302 '/' shift, and go to state 15
303
304 '*' [reduce using rule 8 (exp)]
305 '+' [reduce using rule 8 (exp)]
306 '-' [reduce using rule 8 (exp)]
307 '/' [reduce using rule 8 (exp)]
308 $default reduce using rule 8 (exp)

```

El comentario en la línea 308 (\$default ...) indica que por defecto, ante cualquier otro terminal que no sea uno de los explícitamente listados, la acción a tomar por el analizador será reducir por la regla 8.

Una revisión a la numeración de la gramática, al comienzo del fichero .output nos permite ver cuál es la regla 8:

```
29 Rules:
30 -----

```

```

31 0:  $start -> input $end
32 1:  input -> /* empty */
33 2:  input -> input line
34 3:  line -> '\n'
35 4:  line -> exp '\n'
36 5:  line -> error '\n'
37 6:  exp -> NUM
38 7:  exp -> VAR
39 8:  exp -> VAR '=' exp
40 9:  exp -> exp '+' exp
41 10: exp -> exp '-' exp
42 11: exp -> exp '*' exp
43 12: exp -> exp '/' exp
44 13: exp -> '-' exp
45 14: exp -> '(' exp ')'

```

Efectivamente, es la regla de asignación $\text{exp} \rightarrow \text{VAR} '=' \text{exp}$. El conflicto aparece por que los terminales $* + - /$ están en el conjunto $\text{FOLLOW}(\text{exp})$ y también cabe esperarlos respectivamente en las reglas 9, 10, 11 y 12 ya que el estado 25 contiene:

```

294  exp -> exp . '+' exp  (Rule 9)
295  exp -> exp . '-' exp  (Rule 10)
296  exp -> exp . '*' exp  (Rule 11)
297  exp -> exp . '/' exp  (Rule 12)

```

Estamos ante un caso en el que se aplica el consejo número 8. Los items de la forma $B \rightarrow \beta \uparrow \gamma$, son los de la forma $\text{exp} \rightarrow \text{exp} . '+' \text{exp}$, etc. El item de la forma $A \rightarrow \alpha \uparrow$ es en este caso $\text{exp} \rightarrow \text{VAR} '=' \text{exp}$.

En efecto, en una expresión como $a = 4 + 3$ se produce una ambigüedad. ¿Debe interpretarse como $(a = 4) + 3$? ¿O bien como $a = (4 + 3)$?. La primera interpretación corresponde a reducir por la regla 8. La segunda a desplazar al estado 13. En este ejemplo, el conflicto se resuelve haciendo que tenga prioridad el desplazamiento, dando menor prioridad al terminal $=$ que a los terminales $* + - /$.

Ejercicio 7.17.1. ¿Que ocurre en el ejemplo anterior si dejamos que yacc aplique las reglas por defecto?

7.18. Práctica: Un C simplificado

Escriba un analizador sintáctico usando `Parse::Yapp` para el siguiente lenguaje. La descripción utiliza una notación tipo BNF: las llaves indican 0 o mas repeticiones y los corchetes opcionalidad.

program	→	definitions { definitions }	
definitions	→	datadefinition functiondefinition	
datadefinition	→	basictype declarator { ',' declarator } ';'	
declarator	→	ID { '[' constantexp ']' }	
functiondefinition	→	[basictype] functionheader functionbody	
basictype	→	INT CHAR	
functionheader	→	ID '(' [parameters] ')'	
parameters	→	basictype declarator { ',' basictype declarator }	
functionbody	→	'{ ' { datadefinition } { statement } '}'	
statement	→	[exp] ';' '{ ' { datadefinition } { statement } '}' IF '(' exp ')' statement [ELSE statement] WHILE '(' exp ')' statement RETURN [exp] ';'	Su
constantexp	→	exp	
exp	→	lvalue '=' exp lvalue '+=' exp exp '&&' exp exp ' ' exp exp '==' exp exp '!=' exp exp '<' exp exp '>' exp exp '<=' exp exp '>=' exp exp '+' exp exp '-' exp exp '*' exp exp '/' exp unary	
unary	→	'++' lvalue '--' lvalue primary	
primary	→	'(' exp ')' ID '(' [argumentlist] ') ' lvalue NUM CHARACTER	
lvalue	→	ID { '[' exp ']' }	
argumentlist	→	exp { ',' exp }	

analizador, además de seguir los consejos explicitados en la sección 7.17, deberá cumplir las siguientes especificaciones:

1. Método de Trabajo

Parta de la definición BNF y proceda a introducir las reglas poco a poco:

```

1 %token declarator basictype functionheader functionbody
2 %%
3 program: definitionslist
4         ;
5
6 definitionslist: definitions definitionslist
7                 | definitions
8                 ;
9
10 definitions: datadefinition
11              | functiondefinition
12              ;
13 datadefinition: basictype declaratorlist ';'
14                ;
15
16 declaratorlist: declarator ',' declaratorlist
17                | declarator
18                ;
19 functiondefinition: basictype functionheader functionbody
20                    | functionheader functionbody
21                    ;

```

23 %

- ```
19 functiondefinition: basictype functionheader functionbody
20 | functionheader functionbody
21 ;
```

- ```
:!yapp %
para recompilar:
```

[illegible]

Esto llama a `yapp` con el fichero bajo edición. Si hay errores los detectaré enseguida.

- f) Insisto, procure detectar la aparición de un conflicto lo antes posible. Es terrible tener que limpiar una gramática llena de conflictos.
- g) Ponga nombres significativos a las variables y terminales. Por favor, no los llame `d1`, `d2`, etc.
- h) Cuando esté en el proceso de construcción de la gramática y aún le queden por rellenar variables sintácticas, declárelas como terminales mediante `%token` como en el código que aparece encima. De esta manera evitará las quejas de `yapp`.

2. Resolución de Conflictos

Las operaciones de asignación tienen la prioridad mas baja, seguidas de las lógicas, los test de igualdad y después de los de comparación, a continuación las aditivas, multiplicativas y por último los **unary** y **primary**. Expresa la asociatividad natural y la prioridad especificada usando los mecanismos que **yapp** provee al efecto.

La gramática es ambigua, ya que para una sentencia como

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$$

existen dos árboles posibles: uno que asocia el “else” con el primer “if” y otra que lo asocia con el segundo. Los dos árboles corresponden a las dos posibles parentizaciones:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1 \text{ else } S_2)$$

Esta es la regla de prioridad usada en la mayor parte de los lenguajes: un “else” casa con el “if” mas cercano. La otra posible parentización es:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1) \text{ else } S_2$$

Utilice los mecanismos de priorización proporcionados por **yapp** para resolver el conflicto shift-reduce generado. ¿Es correcta en este caso particular la conducta a la que da lugar la acción **yapp** por defecto?

3. Analizador Léxico

Además del tipo de terminal y su valor el analizador léxico deberá devolver el número de línea. El analizador léxico deberá aceptar comentarios C. En la gramática, el terminal **CHARACTER** se refiere a caracteres entre comillas simples (por ejemplo ‘a’).

Se aconseja que las palabras reservadas del lenguaje no se traten con expresiones regulares específicas sino que se capturen en el patrón de identificador `[a-z_]\w+`. Se mantiene para ello un hash con las palabras reservadas que es inicializado al comienzo del programa. Cuando el analizador léxico encuentra un identificador mira en primer lugar en dicho hash para ver si es una palabra reservada y, si lo es, devuelve el terminal correspondiente. En caso contrario se trata de un identificador.

4. Recuperación de Errores

Extienda la práctica con reglas para la recuperación de errores. Para las listas, siga los consejos dados en la sección 7.2. En aquellos casos en los que la introducción de las reglas de recuperación produzca ambigüedad, resuelva los conflictos.

5. Árbol de Análisis Abstracto

La semántica del lenguaje es similar a la del lenguaje C (por ejemplo, las expresiones lógicas se tratan como expresiones enteras). El analizador deberá producir un árbol sintáctico abstracto. Como se hizo para el lenguaje Tutu introducido en el capítulo 4, cada clase de nodo deberá corresponderse con una clase Perl. Por ejemplo, para una regla como

$$\text{exp '}' \text{ exp}$$

la acción asociada sería algo parecido a

```
{ bless [ $_[1], $_[3]], 'MULT' }
```

donde usamos un array anónimo. Mejor aún es usar un hash anónimo:

```
{ bless { LEFT => $_[1], RIGHT => $_[3]}, 'MULT' }
```

Defina formalmente el árbol especificando la gramática árbol correspondiente a su diseño (repase la sección 4.9.1). Introduzca en esta parte la tabla de símbolos. La tabla de símbolos es, como en el compilador de Tutu, una lista de referencias a hashes conteniendo las tablas de símbolos locales a cada bloque. En cada momento, la lista refleja el anidamiento de bloques actual. Es posible que, en la declaración de funciones, le interese crear un nuevo bloque en el que guardar los parámetros, de manera que las variables globales queden a nivel 0, los parámetros de una función a nivel 1 y las variables locales de la función a nivel 2 o superior.

7.19. La Gramática de yapp / yacc

En esta sección veremos con mas detalle, la sintaxis de `Parse::Yapp`, usando la propia notación `yapp` para describir el lenguaje. Un programa `yapp` consta de tres partes: la cabeza, el cuerpo y la cola. Cada una de las partes va separada de las otras por el símbolo `%%` en una línea aparte.

```
yapp:  head body tail
head:  headsec '%%'
headsec: #empty
        | decls
decls:  decls decl | decl
body:   rulesec '%%'
rulesec: rulesec rules | rules
rules:  IDENT ':' rhss ';'
tail:   /*empty*/
        | TAILCODE
```

7.19.1. La Cabecera

En la cabecera se colocan las declaraciones de variables, terminales, etc.

```
decl:  '\n'
        | TOKEN typedecl symlist '\n'
        | ASSOC typedecl symlist '\n'
        | START ident '\n'
        | HEADCODE '\n'
        | UNION CODE '\n'
        | TYPE typedecl identlist '\n'
        | EXPECT NUMBER '\n'

typedecl: # empty
        | '<' IDENT '>'
```

El terminal `START` se corresponde con una declaración `%start` indicando cual es el símbolo de arranque de la gramática. Por defecto, el símbolo de arranque es el primero de la gramática.

El terminal `ASSOC` está por los terminales que indican precedencia y asociatividad. Esto se ve claro si se analiza el contenido del fichero `YappParse.y` (??) en el que se puede encontrar el código del analizador léxico del módulo `Parse::Yapp`. El código dice:

```
...
if($lexlevel == 0) {# In head section
    $$input=~/\G%(left|right|nonassoc)/gc
```



```

and return('ASSOC',[ uc($1), $lineno[0] ]);
    $$input=~/\G%(start)/gc
and return('START',[ undef, $lineno[0] ]);
    $$input=~/\G%(expect)/gc
and return('EXPECT',[ undef, $lineno[0] ]);
    $$input=~/\G%{/gc
...

```

La variable `$lexlevel` indica en que sección nos encontramos: cabecera, cuerpo o cola. El terminal `EXPECT` indica la presencia de una declaración `%expect` en el fuente, la cual cuando es seguida de un número indica el numero de conflictos shift-reduce que cabe esperar. Use `EXPECT` si quiere silenciar las advertencias de `yapp` sobre la presencia de conflictos cuya resolución automática considere correcta.

7.19.2. La Cabecera: Diferencias entre `yacc` y `yapp`

Las declaraciones de tipo correspondientes a `%union` y a las especificaciones de tipo entre símbolos menor mayor (`<tipo>`) en declaraciones `token` y `%type` no son usadas por `yapp`. Estas declaraciones son necesarias cuando el código de las acciones semánticas se escribe en C como es el caso de `yacc` y `bison`. Sigue un ejemplo de programa `yacc/bison` que usa declaraciones `%union` y de tipo para los atributos:

```

1 %{
2 #include <stdio.h>
3
4 #define CLASE(x) ((x == 1)?"global":"local")
5 #define TIPO(x) ((x == 1)?"float":"integer")
6 %}
7
8 %union {
9     int n;      /* enumerado */
10    char *s;     /* cadena */
11 }
12
13 %token <n> FLOAT INTEGER
14 %token <n> GLOBAL
15 %token <n> LOCAL
16 %token <s> NAME
17 %type <n> class type
18
19 %%

```

La declaración `%union` de la línea 8 indica que los atributos son de dos tipos: enteros y punteros a caracteres. El nombre del campo es posteriormente usado en las declaraciones de las líneas 13-17 para indicar el tipo del atributo asociado con la variable o con el terminal. Así, la declaración de la línea 13 indica que los terminales `FLOAT` e `INTEGER` son de tipo entero, mientras que la declaración de la línea 16 nos dice que el terminal `NAME` es de tipo cadena.

```

29 class
30 : GLOBAL { $$ = 1; }
31 | LOCAL { $$ = 2; }
32 ;
33
34 type
35 : FLOAT { $$ = 1; }
36 | INTEGER { $$ = 2; }
37 ;

```

La información proveída sobre los tipos permite a `yacc` introducir automáticamente en el código C producido los *typecasting* o ahormados para las asignaciones de las líneas 30-31 y 35-36. Obsérvese que en `yacc` el atributo de la variable en la parte izquierda se denota por `$$`.

Otra diferencia entre `yacc` y `yapp` es que en `yacc` los atributos de la parte derecha no constituyen un vector, denotándose por `$1`, `$2`, `$3` ...

En ocasiones `yacc` no puede determinar el tipo de un atributo. En particular cuando se habla del atributo asociado con una acción intermedia, ya que esta no tiene variable sintáctica asociada explícitamente o bien cuando se habla de los atributos de símbolos que están a la izquierda de la reducción actual (véase la sección 8.28). Los atributos de símbolos a la izquierda de la producción actual se denotan en `yacc` por números no positivos `$0`, `$-1`, `$-2`

En estos casos el programador deberá especificar explícitamente el tipo del atributo usando la notación `$<tipo>#`. Donde `tipo` es uno de los campos de la `union` y `#` es el numeral del símbolo correspondiente:

```
39 namelist
40   : NAME    { printf("%s de clase %s, tipo %s\n", $1, CLASE($<n>-1), TIPO($<n>0)); }
41   | namelist ',' NAME
42       { printf("%s de clase %s, tipo %s\n", $3, CLASE($<n>-1), TIPO($<n>0)); }
43   ;
44 %%
```

7.19.3. El Cuerpo

El cuerpo de un programa `yapp` contiene la gramática y las acciones

```
rhss:   rhss '|' rule | rule
rule:   rhs prec epscode | rhs
rhs:    #empty
        | rhselts
rhselts: rhselts rhselt | rhselt
rhselt: symbol | code
prec:   PREC symbol
epscode: # vacio
        | code
code:   CODE
```

Las acciones semánticas (variable sintáctica `code` y terminal `CODE`) se ejecutan siempre que ocurre una reducción por una regla y, en general, devuelven un valor semántico. El código de la acción se copia verbatim en el analizador. La estrategia usada por el analizador léxico es contar las llaves abrir y cerrar en el texto. Véase el correspondiente fragmento del analizador léxico:

```
....
$lineno[0]=$lineno[1];
....
$$input=~/\G{/gc
and do {
    my($level,$from,$code);

    $from=pos($$input);
    $level=1;
    while($$input=~/([{])/gc) {
        substr($$input,pos($$input)-1,1) eq '\{' #Quoted
        and next;
        $level += ($1 eq '{' ? 1 : -1) or last;
    }
}
```

```

    $level and _SyntaxError(2,"Unmatched { opened line $lineno[0]",-1);
    $code = substr($$input,$from,pos($$input)-$from-1);
    $lineno[1]+= $code=~tr/\n//;
    return('CODE',[ $code, $lineno[0] ]);
};

```

Las llaves dentro de cadenas y comentarios no son significativas en la cuenta. El problema es que el reconocimiento de cadenas en Perl es mas difícil que en otros lenguajes: existe toda una variedad de formas de denotar una cadena. Por tanto, si el programador usuario de `yapp` necesita escribir una llave dentro de una cadena de doble comilla, deberá escaparla. Si la cadena es de simple comilla escaparla no es solución, pues aparecería el símbolo de escape en la cadena. En ese caso se deberá añadir un comentario con la correspondiente falsa llave. Siguen algunos ejemplos tomados de la documentación de `Parse::Yapp`

```

"{ My string block }"
"\{ My other string block \}"
qq/ My unmatched brace \} /

# Casamos con el siguiente: {
q/ for my closing brace } / #

q/ My opening brace { /
# debe cerrarse: }

```

Ejercicio 7.19.1. *Genere programas de prueba `yapp` con cadenas que produzcan confusión en el analizador y observe el comportamiento. Pruébelas en las diferentes secciones en las que puede ocurrir código: en la cabecera, en el cuerpo y en la cola.*

7.19.4. La Cola: Diferencias entre `yacc` y `yapp`

La cola de un program `yapp` contiene las rutinas de soporte.

```

tail:  /*empty*/
      |  TAILCODE

```

el terminal `TAILCODE` al igual que los terminales `CODE` y `HEADCODE` indican que en ese punto se puede encontrar código Perl. La detección de `TAILCODE` y `HEADCODE` son mas sencillas que las de `CODE`.

La cola de un programa `yacc` es similar. Para el programa `yacc` cuya cabecera y cuerpo se mostraron en la sección 7.19.2 la cola es:

```

1  %%
2
3  extern FILE * yyin;
4
5  main(int argc, char **argv) {
6      if (argc > 1) yyin = fopen(argv[1],"r");
7      /* yydebug = 1;
8      */
9      yyparse();
10 }
11
12 yyerror(char *s) {
13     printf("%s\n",s);
14 }

```

La declaración del manejador de fichero `yyin` en la línea 14 referencia el archivo de entrada para el analizador. La variable (comentada, línea 7) `yydebug` controla la información para la depuración de la gramática. Para que sea realmente efectiva, el programa deberá además compilarse definiendo la macro `YYDEBUG`. Sigue un ejemplo de `Makefile`:

```
1  inherited: y.tab.c lex.yy.c
2      gcc -DYYDEBUG=1 -g -o inherited1 y.tab.c lex.yy.c
3  y.tab.c y.tab.h: inherited1.y
4      yacc -d -v inherited1.y
5  lex.yy.c: inherited1.l y.tab.h
6      flex -l inherited1.l
7  clean:
8      - rm -f y.tab.c lex.yy.c *.o core inherited1
```

Al compilar tenemos:

```
pl@nereida:~/src/inherited$ make
yacc -d -v inherited1.y
flex -l inherited1.l
gcc -DYYDEBUG=1 -g -o inherited1 y.tab.c lex.yy.c
pl@nereida:~/src/inherited$ ls -ltr
total 232
-rw-r----- 1 pl users 242 Dec 10 2003 Makefile
-rw-r----- 1 pl users 404 Dec 10 2003 inherited1.l
-rw-r----- 1 pl users 878 Dec 10 2003 inherited1.y
-rw-rw---- 1 pl users 1891 Jan 26 15:41 y.tab.h
-rw-rw---- 1 pl users 30930 Jan 26 15:41 y.tab.c
-rw-rw---- 1 pl users 2365 Jan 26 15:41 y.output
-rw-rw---- 1 pl users 44909 Jan 26 15:41 lex.yy.c
-rwxrwx--x 1 pl users 56336 Jan 26 15:41 inherited1
```

7.19.5. El Análisis Léxico en yacc: flex

El analizador léxico para yacc desarrollado en las secciones anteriores ha sido escrito usando la variante `flex` del lenguaje LEX. Un programa `flex` tiene una estructura similar a la de un program yacc con tres partes: cabeza, cuerpo y cola separados por `%%`. Veamos como ejemplo de manejo de `flex`, los contenidos del fichero `flex inherited1.l` utilizado en las secciones anteriores:

```
1  %{
2  #include <string.h>
3  #include "y.tab.h"
4  %}
5  id [A-Za-z_][A-Za-z_0-9]*
6  white [ \t\n]+
7  %%
8  global { return GLOBAL; }
9  local { return LOCAL; }
10 float { return FLOAT; }
11 int { return INTEGER; }
12 {id} { yylval.s = strdup(yytext); return NAME; }
13 {white} { ; }
14 , { return yytext[0]; }
15 . { fprintf(stderr,"Error. carácter inesperado.\n"); }
16 %%
17 int yywrap() { return 1; }
```

La cabeza contiene declaraciones C así como definiciones regulares. El fichero `y.tab.h` que es incluido en la línea 3, fué generado por `yacc` y contiene, entre otras cosas, la información recolectada por `yacc` sobre los tipos de los atributos (declaración `%union`) y la enumeración de los terminales. Es, por tanto, necesario que la compilación con `yacc` preceda a la compilación con `flex`. La información en `y.tab.h` es usada por el analizador léxico para “sincronizarse” con el analizador sintáctico. Se definen en las líneas 5 y 6 las macros para el reconocimiento de identificadores (`id`) y blancos (`white`). Estas macros son llamadas en el cuerpo en las líneas 12 y 13. La estructura del cuerpo consiste en parejas formadas por una definición regular seguidas de una acción. La variable `yylval` contiene el atributo asociado con el terminal actual. Puesto que el token `NAME` fué declarado del tipo cadena (véase 7.19.2), se usa el correspondiente nombre de campo `yylval.s`. La cadena que acaba de casar queda guardada en la variable `yytext`, y su longitud queda en la variable entera global `yylen`.

Una vez compilado con `flex` el fuente, obtenemos un fichero denominado `lex.yy.c`. Este fichero contiene la rutina `yylex()` que realiza el análisis léxico del lenguaje descrito.

La función `yylex()` analiza las entradas, buscando la secuencia mas larga que casa con alguna de las expresiones regulares y ejecuta la correspondiente acción. Si no se encuentra ningun emparejamiento se ejecuta la regla “por defecto”, que es:

```
(.|\n) { printf("%s",yytext); }
```

Si encuentran dos expresiones regulares con las que la cadena mas larga casa, elige la que figura primera en el programa `flex`.

Una vez que se ha ejecutado la correspondiente acción, `yylex()` continúa con el resto de la entrada, buscando por subsiguientes emparejamientos. Así continúa hasta encontrar un final de fichero en cuyo caso termina, retornando un cero o bien hasta que una de las acciones explícitamente ejecuta una sentencia `return`.

Cuando el analizador léxico alcanza el final del fichero, el comportamiento en las subsiguientes llamadas a `yylex` resulta indefinido. En el momento en que `yylex` alcanza el final del fichero llama a la función `yywrap`, la cual retorna un valor de 0 o 1 según haya mas entrada o no. Si el valor es 0, la función `yylex` asume que la propia `yywrap` se ha encargado de abrir el nuevo fichero y asignárselo a `yyin`.

7.19.6. Práctica: Uso de Yacc y Lex

Use `yacc` y `flex` para completar los analizadores sintáctico y léxico descritos en las secciones 7.19.2, 7.19.4 y 7.19.5. La gramática en cuestión es similar a la descrita en la sección 8.28. Usando la variable `yydebug` y la macro `YYDEBUG` analice el comportamiento para la entrada `global float x,y`.

7.20. El Analizador Ascendente Parse::Yapp

El program `yapp` es un traductor y, por tanto, constituye un ejemplo de como escribir un traductor. El lenguaje fuente es el lenguaje `yacc` y el lenguaje objeto es `Perl`. Como es habitual en muchos lenguajes, el lenguaje objeto se ve `.expandido` con un conjunto de funciones de soporte. En el caso de `yapp` estas funciones de soporte, son en realidad métodos y están en el módulo `Parse::Yapp::Driver`. Cualquier módulo generado por `yapp` hereda de dicho módulo (véase por ejemplo, el módulo generado para nuestro ejemplo de la calculadora, en la sección 7.4).

Como se ve en la figura 7.3, los módulos generados por `yapp` heredan y usan la clase `Parse::Yapp::Driver` la cual contiene el analizador sintáctico LR genérico. Este módulo contiene los métodos de soporte visibles al usuario `YYParse`, `YYData`, `YYError`, `YYSemval`, etc.

La figura 7.3 muestra además el resto de los módulos que conforman el “compilador” `Parse::Yapp`. La herencia se ha representado mediante flechas continuas. Las flechas punteadas indican una relación de uso entre los módulos. El guión `yapp` es un programa aparte que es usado para producir el correspondiente módulo desde el fichero conteniendo la gramática.

(Para ver el contenido de los módulos, descarge `yapp` desde CPAN:

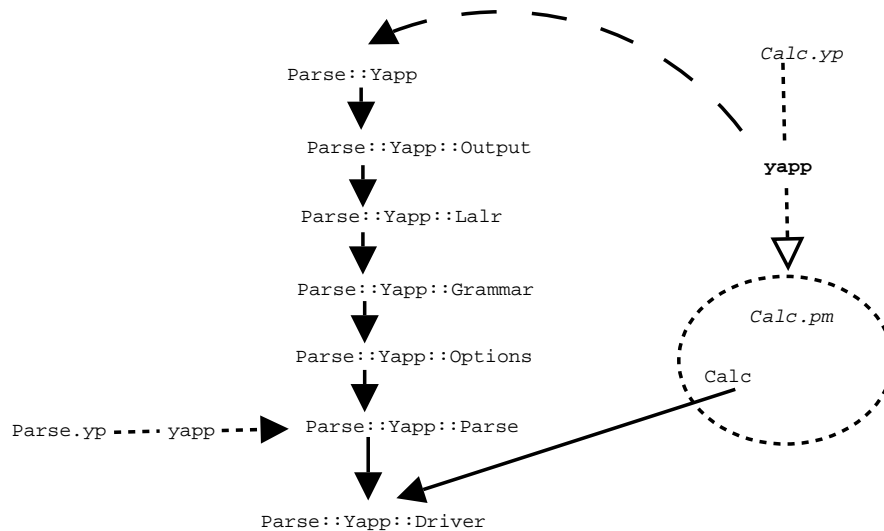


Figura 7.3: Esquema de herencia de `Parse::Yapp`. Las flechas continuas indican herencia, las punteadas uso. La clase `Calc` es implementada en el módulo generado por `yapp`

`\http://search.cpan.org/~fdesar/Parse-Yapp-1.05/lib/Parse/Yapp.p`

o bien desde uno de nuestros servidores locales; en el mismo directorio en que se guarda la versión HTML de estos apuntes encontrará una copia de `Parse-Yapp-1.05.tar.gz`). La versión a la que se refiere este capítulo es la 1.05.

El módulo `Parse/Yapp/Yapp.pm` se limita a contener la documentación y descansa toda la tarea de análisis en los otros módulos. El módulo `Parse/Yapp/Output.pm` contiene los métodos `_CopyDriver` y `Output` los cuales se encargan de escribir el analizador: partiendo de un esqueleto genérico rellenan las partes específicas a partir de la información computada por los otros módulos.

El módulo `Parse/Yapp/Options.pm` analiza las opciones de entrada. El módulo `Parse/Yapp/Lalr.pm` calcula las tablas de análisis LALR. Por último el módulo `Parse/Yapp/Grammar` contiene varios métodos de soporte para el tratamiento de la gramática.

El módulo `Parse::Yapp::Driver` contiene el método `YYparse` encargado del análisis. En realidad, el método `YYparse` delega en el método privado `_Parse` la tarea de análisis. Esta es la estructura del analizador genérico usado por `yapp`. Léalo con cuidado y compare con la estructura explicada en la sección 8.25.

```

1 sub _Parse {
2   my($self)=shift;
3
4   my($rules,$states,$lex,$error)
5     = @$self{ 'RULES', 'STATES', 'LEX', 'ERROR' };
6   my($errstatus,$nberror,$token,$value,$stack,$check,$dotpos)
7     = @$self{ 'ERRST', 'NBERR', 'TOKEN', 'VALUE', 'STACK', 'CHECK', 'DOTPOS' };
8
9   $$errstatus=0;
10  $$nberror=0;
11  ($$token,$$value)=(undef,undef);
12  @$stack=( [ 0, undef ] ); # push estado 0
13  $$check='';

```

La componente 0 de `@$stack` es el estado, la componente 1 es el atributo.

```

14
15  while(1) {
16    my($actions,$act,$stateno);

```

```

17
18     $stateno=$$stack[-1][0];      # sacar el estado en el top de
19     $actions=$$states[$stateno]; # la pila

```

\$states es una referencia a un vector. Cada entrada \$\$states[\$stateno] es una referencia a un hash que contiene dos claves. La clave ACTIONS contiene las acciones para ese estado. La clave GOTOS contiene los saltos correspondientes a ese estado.

```

20
21     if (exists($$actions{ACTIONS})) {
22         defined($$token) or do {
23             ($$token,$$value)=&$lex($self); # leer siguiente token
24         };
25
26         # guardar en $act la acción asociada con el estado y el token
27         $act = exists($$actions{ACTIONS}{$$token})?
28             $$actions{ACTIONS}{$$token} :
29             exists($$actions{DEFAULT})? $$actions{DEFAULT} : undef;
30     }
31     else { $act=$$actions{DEFAULT}; }

```

La entrada DEFAULT de una acción contiene la acción a ejecutar por defecto.

```

32
33     defined($act) and do {
34         $act > 0 and do { # $act >0 indica shift
35             $$errstatus and do { --$$errstatus; };

```

La línea 35 esta relacionada con la recuperación de errores. Cuando yapp ha podido desplazar varios terminales sin que se produzca error considerará que se ha recuperado con éxito del último error.

```

36         # Transitar: guardar (estado, valor)
37         push(@$stack,[ $act, $$value ]);
38         $$token ne ''      #Don't eat the eof
39         and $$token=$$value=undef;
40         next; # siguiente iteración
41     };

```

A menos que se trate del final de fichero, se reinicializa la pareja (\$\$token, \$\$value) y se repite el bucle de análisis. Si \$act es negativo se trata de una reducción y la entrada \$\$rules[-\$act] es una referencia a un vector con tres elementos: la variable sintáctica, la longitud de la parte derecha y el código asociado:

```

43         # $act < 0, indica reduce
44         my($lhs,$len,$code,@sempar,$semval);
45
46         #obtenemos una referencia a la variable,
47         #longitud de la parte derecha, referencia
48         #a la acción
49         ($lhs,$len,$code)=@{$$rules[-$act]};
50         $act or $self->YYAccept();

```

Si \$act es cero indica una acción de aceptación. El método YYAccept se encuentra en Driver.pm. Simplemente contiene:

```

sub YYAccept {
    my($self)=shift;

```

```

        ${$self{CHECK}}='ACCEPT';
        undef;
    }

```

Esta entrada será comprobada al final de la iteración para comprobar la condición de aceptación (a través de la variable `$check`, la cuál es una referencia).

```

51     $$dotpos=$len; # dotpos es la longitud de la regla
52     unpack('A1',$lhs) eq '@'      #In line rule
53     and do {
54         $lhs =~ /\^@([0-9]+\-[0-9]+)$/
55         or die "In line rule name '$lhs' ill formed: ".
56             "report it as a BUG.\n";
57         $$dotpos = $1;
58     };

```

En la línea 52 obtenemos el primer carácter en el nombre de la variable. Las acciones intermedias en `yapp` producen una variable auxiliar que comienza por `@` y casa con el patrón especificado en la línea 54. Obsérvese que el número después del guión contiene la posición relativa en la regla de la acción intermedia.

```

60     @sempar = $$dotpos ?
61         map { $$_[1] } @$stack[ -$$dotpos .. -1 ] : ();

```

El array `@sempar` se inicia a la lista vacía si `$len` es nulo. En caso contrario contiene la lista de los atributos de los últimos `$$dotpos` elementos referenciados en la pila. Si la regla es intermedia estamos haciendo referencia a los atributos de los símbolos a su izquierda.

```

62     $semval = $code ? &$code( $self, @sempar ) :
63         @sempar ? $sempar[0] : undef;

```

Es en este punto que ocurre la ejecución de la acción. La subrutina referenciada por `$code` es llamada con primer argumento la referencia al objeto analizador `$self` y como argumentos los atributos que se han computado previamente en `@sempar`. Si no existe tal código se devuelve el atributo del primer elemento, si es que existe un tal primer elemento.

El valor retornado por la subrutina/acción asociada es guardado en `$semval`.

```

65     splice(@$stack,-$len,$len);

```

La función `splice` toma en general cuatro argumentos: el array a modificar, el índice en el cual es modificado, el número de elementos a suprimir y la lista de elementos extra a insertar. Aquí, la llamada a `splice` cambia los elementos de `@$stack` a partir del índice `-$len`. El número de elementos a suprimir es `$len`. A continuación se comprueba si hay que terminar, bien porque se ha llegado al estado de aceptación (`$$check eq 'ACCEPT'`) o porque ha habido un error fatal:

```

    $$check eq 'ACCEPT' and do { return($semval); };
    $$check eq 'ABORT' and do { return(undef); };

```

Si las cosas van bien, se empuja en la cima de la pila el estado resultante de transitar desde el estado en la cima con la variable sintáctica en el lado izquierdo:

```

    $$check eq 'ERROR' or do {
        push(@$stack, [ $$states[$$stack[-1][0]]{GOTOS}{$lhs}, $semval ]);
        $$check='';
        next;
    };

```


La expresión `$$states[$$stack[-1][0]` es una referencia a un hash cuya clave `GOTOS` contiene una referencia a un hash conteniendo la tabla de transiciones del estado en la cima de la pila (`$stack[-1][0]`). La entrada de clave `$lhs` contiene el estado al que se transita al ver la variable sintáctica de la izquierda de la regla de producción. El atributo asociado es el devuelto por la acción: `$semval`.

```

        $$check='';

    }; # fin de defined($act)

    # Manejo de errores: código suprimido
    ...

}
}#_Parse

```

... y el bucle `while(1)` de la línea 15 continúa. Compare este código con el pseudo-código introducido en la sección 8.25.

7.21. La Estructura de Datos Generada por `YappParse.yp`

El fichero `YappParse.yp` contiene la gramática `yapp` del lenguaje `yacc`¹. Además de las dos rutinas de soporte típicas, la de tratamiento de errores `_Error` y la de análisis léxico `_Lexer`, el fichero contiene una subrutina para el manejo de las reglas `_AddRules` y otra rutina `Parse` la cuál actúa como *wrapper* o filtro sobre el analizador `YYParse`.

Durante el análisis sintáctico de un programa `yapp` se construye una estructura de datos para la posterior manipulación y tratamiento de la gramática. Como ejemplo usaremos la gramática:

```

pl@nereida:~/src/perl/Parse-AutoTree/trunk/scripts$ cat -n int.yp
1  %right '+'
2  %left 'a'
3  %nonassoc 'b'
4  %%
5  S:  /* empty rule */                { print "S -> epsilon\n" }
6  |   'a' { print "Intermediate\n"; } S 'b' { print "S -> a S b\n" }
7  |   '+' S '+' %prec 'a'             { print "S -> + S + prec a\n" }
8  ;
9  %%
10
11 sub _Error {
12     exists $_[0]->YYData->{ERRMSG}
13     and do {
14         print $_[0]->YYData->{ERRMSG};
15         delete $_[0]->YYData->{ERRMSG};
16         return;
17     };
18     print "Syntax error.\n";
19 }
20
21 sub _Lexer {
22     my($parser)=shift;
23
24     defined($parser->YYData->{INPUT})

```

¹La versión a la que se refiere esta sección es la 1.05 (`Parse-Yapp-1.05.tar.gz`)

```

25     or $parser->YYData->{INPUT} = <STDIN>
26     or return('',undef);
27
28     $parser->YYData->{INPUT}=~s/^[ \t\n]//;
29
30     for ($parser->YYData->{INPUT}) {
31         s/^(.)//s and return($1,$1);
32     }
33 }
34
35 sub Run {
36     my($self)=shift;
37     $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error, yydebug => 0x1F );
38 }

```

Para construir la estructura podemos usar la siguiente subrutina:

```

sub Parse {
    my $grammar = shift;

    my $x = new Parse::Yapp::Parse;
    my $r = $x->Parse($grammar);

    return $r;
}

```

La llamada a Parse produce la siguiente estructura de datos:

```

nereida:~/src/perl/Parse-AutoTree/trunk/scripts> grammar.pl int.yp
$VAR1 = {
    'START' => 'S', # Símbolo de arranque
    'SYMS' => { 'S' => 5, 'b' => 3, 'a' => 2, '+' => 1 }, # Símbolo => línea
    'TERM' => {
        'b' => [ 'NONASSOC', 2 ], # terminal => [ Asociatividad, precedencia ]
        'a' => [ 'LEFT', 1 ], # terminal => [ Asociatividad, precedencia ]
        '+' => [ 'RIGHT', 0 ] }, # terminal => [ Asociatividad, precedencia ]
    # Si el terminal no tiene precedencia toma la forma terminal => []
    'NTERM' => { 'S' => [ '1', '3', '4' ] }, # variable => [ indice en RULES de las reglas de S
    'PREC' => { 'a' => 1 }, # Terminales que son usados en una directiva %prec
    'NULL' => { 'S' => 1 }, # Variables que producen vacío
    'EXPECT' => 0, # Número de conflictos esperados
    'RULES' => [
        [ '$start', [ 'S', '' ], undef, undef ], # Regla de superarranque
        [
            'S', [], # producción
            undef, # precedencia explícita de la regla
            [ ' print "S -> epsilon\n" ', 5 ] # [ acción asociada, línea ]
        ],
        [
            '@1-1', [], # Regla intermedia: Variable temporal
            undef,
            [ ' print "Intermediate\n"; ', 6 ]
        ],
        [
            'S', [ 'a', '@1-1', 'S', 'b' ],

```

```

    undef,
    [ ' print "S -> a S b\n" ', 6 ]
],
[
    'S', [ '+', 'S', '+' ],
    1, # precedencia explícita de la regla
    [ ' print "S -> + S + prec a\n" ', 7 ]
]
],
'HEAD' => undef, # Código de cabecera
'TAIL' => [ '... código de cola ...', 9 ], # Línea en la que comienza la sección de cola
};

```

Las componentes del hash que aparece arriba se corresponden con diversas variables usadas por `YYParse` durante el análisis. La correspondencia se establece dentro del método `Parse` cuando se hace la asignación:

```

@$parsed{ 'HEAD', 'TAIL', 'RULES', 'NTERM', 'TERM',
          'NULL', 'PREC', 'SYMS', 'START', 'EXPECT' }
=
    ( $head, $tail, $rules, $nterm, $term,
      $nullable, $precterm, $syms, $start, $expect);

```

esta asignación es la que crea el hash. Las variables con identificadores en minúsculas son usadas en el analizador. Son visibles en todo el fichero ya que, aunque declaradas léxicas, su declaración se encuentra en la cabecera del analizador:

```

%{
require 5.004;

use Carp;

my($input,$lexlevel,@lineno,$nberr,$prec,$labelno);
my($syms,$head,$tail,$token,$term,$nterm,$rules,$precterm,$start,$nullable);
my($expect);

%}

```

7.22. Práctica: El Análisis de las Acciones

Modifique el código de `YappParse.y` para que el análisis léxico de las secciones de código (`HEADCODE`, `CODE` y `TAILCODE`) se haga a través de las correspondientes rutinas proveída como parámetros para el análisis por el usuario. La idea es ofrecer un primer paso que facilite la generación de analizadores en diferentes lenguajes Perl, C, etc.

Estudie el módulo `Text::Balanced`. Basándose en las funciones

`extract_codeblock` y `extract_quotelike`

del módulo `Text::Balanced`, resuelva el problema del reconocimiento de código Perl dentro del analizador léxico de `Parse::Yapp`, evitando forzar al usuario en la escritura de “llaves fantasma”. Compare el rendimiento de esta solución con la que provee `Yapp`. Para analizar el *rendimiento* use el módulo `Benchmark`.

¿Cuáles son sus conclusiones? ¿Qué es mejor?

7.23. Práctica: Autoacciones

Extienda `Parse::Yapp` con una directiva `%autoaction CODE` la cuál cambia la acción por defecto. Cuando una regla de producción no tenga una acción asociada, en vez de ejecutarse la acción `yapp`

por defecto se ejecutará el código especificado en CODE. La directiva podrá aparecer en la parte de cabecera o en el cuerpo del programa yapp en una sólo línea aparte. Si aparece en el cuerpo no debe hacerlo en medio de una regla.

Sigue un ejemplo de uso:

```
%{
use Data::Dumper;
my %tree_name = ('=' => 'eq', '+' => 'plus', '-' => 'minus',
                 '*' => 'times', '/' => 'divide');
%}
%right  '='
%left  '-' '+'
%left  '*' '/'
%left  NEG
%autoaction { [$tree_name{$_[2]}, $_[1], $_[3]] }

%%
input:          { undef }
      | input line { undef }
;

line:          '\n'          { undef }
      | exp '\n'          { print Dumper($_[1]); }
      | error '\n'        { $_[0]->YYError }
;

exp:          NUM            { $_[1] }
      | VAR                { $_[1] }
      | VAR '=' exp | exp '+' exp | exp '-' exp | exp '*' exp | exp '/' exp
      | '-' exp %prec NEG   { ['neg', $_[2]] }
      | '(' exp ')'         { $_[2] }
;

%%
```

y un ejemplo de ejecución:

```
$ ./useautoaction1.pl
2+3*4
^D
$VAR1 = [
    'plus',
    '2',
    [ 'times', '3', '4' ]
];
```

Analice la adecuación de los mensajes de error emitidos por el compilador de Perl cuando el código en la auto-acción contiene errores. ¿Son apropiados los números de línea?

Tenga en cuenta los siguientes consejos:

- Cuando compile con yapp su módulo use una orden como: `yapp -m Parse::Yapp::Parse Parse.ypp`. Este es un caso en que el nombre del fichero de salida (`Parse.pm`) y el nombre del package `Parse::Yapp::Parse` no coinciden. Este es un caso en que el nombre del fichero de salida (`Parse.pm`) y el nombre del package `Parse::Yapp::Parse` no coinciden.

- Ahora tiene dos versiones de `Parse::Yapp` en su ordenador. El compilador de Perl va a intentar cargar la instalada. Para ello en su versión del script `yapp` puede incluir una línea que le indique al compilador que debe buscar primero en el lugar en el que se encuentra nuestra librería:

```
BEGIN { unshift @INC, '/home/lhp/Lperl/src/yapp/Parse-Yapp-Auto/lib/' }
```

- ¿Qué estrategia a seguir? Una posibilidad es “hacerle creer” al resto de los módulos en `Yapp` que el usuario ha escrito el código de la autoacción en aquellas reglas en las que no existe código explícito asociado. Es posible realizar esta práctica modificando sólo el fichero `YappParse.y`. El código original `Yapp` usa `undef` para indicar, en el campo adecuado, que una acción no fué definida. La idea es sustituir ese `undef` por el código asociado con la autoacción:

```
my($code)= $autoaction? $autoaction:undef;
```

7.24. Práctica: Nuevos Métodos

Continuemos extendiendo `Yapp`.

- Introduzca en el módulo `Driver.pm` de `Yapp` un método `YYLhs` que devuelva el identificador de la variable sintáctica en el lado izquierdo de la regla de producción por la que se está reduciendo.
- Para tener disponible el lado izquierdo deberá modificar la conducta del analizador LALR (subrutina `_Parse`) para que guarde como un atributo el identificador de dicho noterminal.
- ¿Que identificador se devuelve asociado con las acciones intermedias?

Sigue un ejemplo de como programar haciendo uso de esta y la anterior extensión:

```
%right  '='
%left   '- ' '+'
%left   '*' '/'
%left   NEG
%autoaction { my $n = $#_; bless [@_[1..$n]], $_[0]->YYLhs }
%%
input:
    |   input line
;
line:
    '\n'      { }
    | exp '\n' { [ $_[1] ] }
    | error '\n' { }
;
exp:
    NUM      |   VAR      |   VAR '=' exp
    | exp '+' exp |   exp '-' exp |   exp '*' exp |   exp '/' exp
    | '-' exp %prec NEG
    | '(' exp ')' { [ $_[2] ] }
;
%%
...
```

Veamos la ejecución correspondiente al ejemplo anterior:

```
$ ./uselhs2.pl
2+3*4
$VAR1 = bless(
[
    bless( [], 'input' ),
```

```
[
  bless( [
    bless( [ '2' ], 'exp' ),
    '+',
    bless( [
      bless( [ '3' ], 'exp' ), '*', bless( [ '4' ], 'exp' ) ], 'exp' )
    ], 'exp' )
]
], 'input' );
```

7.25. Práctica: Generación Automática de Árboles

Partiendo de la práctica anterior, introduzca una directiva `%autotree` que de lugar a la construcción del árbol de análisis concreto. La acción de construcción del árbol:

```
{ my $n = $#_; bless [@_[1..$n]], $_[0]->YYLhs }
```

se ejecutará para cualquier regla que no tenga una acción explícita asociada.

7.26. Recuperacion de Errores: Visión Detallada

La subrutina `_Parse` contiene el algoritmo de análisis LR genérico. En esta sección nos concentraremos en la forma en la que se ha implantado en `yapp` la recuperación de errores.

```
1 sub _Parse {
2   my($self)=shift;
3   ...
4   $$errstatus=0; $$nerror=0;
```

La variable `$$errstatus` nos indica la situación con respecto a la recuperación de errores. La variable `$$nerror` contiene el número total de errores.

```
5   ($$token,$$value)=(undef,undef);
6   @$$stack=( [ 0, undef ] ); $$check='';
7   while(1) {
8     my($actions,$act,$stateno);
9     $stateno=$$stack[-1][0];
10    $actions=$$states[$stateno];
11
12    if (exists($$actions{ACTIONS})) {
13      defined($$token) or do { ($$token,$$value)=$$lex($self); };
14      ...
15    }
16    else { $act=$$actions{DEFAULT}; }
```

Si `$act` no esta definida es que ha ocurrido un error. En tal caso no se entra a estudiar si la acción es de desplazamiento o reducción.

```
17    defined($act) and do {
18      $act > 0 and do { #shift
19        $$errstatus and do { --$$errstatus; };
20        ...
21        next;
22      };
```

```

23     #reduce
24     ....
25     $$check eq 'ERROR' or do {
26         push(@$stack, [ $$states[$$stack[-1][0]]{GOTOS}{lhs}, $semval ]);
27         $$check='';
28         next;
29     };
30     $$check='';
31 };

```

Si `$$errstatus` es cero es que estamos ante un nuevo error:

```

32     #Error
33     $$errstatus or do {
34         $$errstatus = 1;
35         &$error($self);
36         $$errstatus # if 0, then YYErrork has been called
37         or next; # so continue parsing
38         ++$$nberror;
39     };

```

Como el error es “nuevo” se llama a la subrutina de tratamiento de errores `&$error`. Obsérvese que no se volverá a llamar a la rutina de manejo de errores hasta que `$$errstatus` vuelva a alcanzar el valor cero. Puesto que `&$error` ha sido escrita por el usuario, es posible que este haya llamado al método `YYErrork`. Si ese es el caso, es que el programador prefiere que el análisis continúe como si la recuperación de errores se hubiera completado.

Ahora se pone `$$errstatus` a 3:

```

47     $$errstatus=3;

```

Cada vez que se logre un desplazamiento con éxito `$$errstatus` será decrementado (línea 19).

A continuación se retiran estados de la pila hasta que se encuentre alguno que pueda transitar ante el terminale especial `error`:

```

48     while(@$stack
49         and (not exists($$states[$$stack[-1][0]]{ACTIONS})
50         or not exists($$states[$$stack[-1][0]]{ACTIONS}{error})
51         or $$states[$$stack[-1][0]]{ACTIONS}{error} <= 0)) {
52         pop(@$stack);
53     }
54     @$stack or do {
55         return(undef);
56     };

```

Si la pila quedó vacía se devuelve `undef`. En caso contrario es que el programador escribió alguna regla para la recuperación de errores. En ese caso, se transita al estado correspondiente:

```

57     #shift the error token
58     push(@$stack, [ $$states[$$stack[-1][0]]{ACTIONS}{error}, undef ]);
59 }
60 #never reached
61 croak("Error in driver logic. Please, report it as a BUG");
62 }#_Parse

```

Un poco antes tenemos el siguiente código:

```

41     $$errstatus == 3  #The next token is not valid: discard it
42     and do {
43         $$token eq ''  # End of input: no hope
44         and do { return(undef); };
45         $$token=$$value=undef;
46     };

```

Si hemos alcanzado el final de la entrada en una situación de error se abandona devolviendo `undef`.

Ejercicio 7.26.1. *Explique la razón para el comentario de la línea 41. Si `$$errstatus` es 3, el último terminal no ha producido un desplazamiento correcto. ¿Porqué?*

A continuación aparecen los códigos de los métodos implicados en la recuperación de errores:

```

sub YYError {
    my($self)=shift;

    ${$self{ERRST}}=0;
    undef;
}

```

El método `YYError` cambia el valor referenciado por `$errstatus`. De esta forma se le da al programador `yapp` la oportunidad de anunciar que es muy probable que la fase de recuperación de errores se haya completado.

Los dos siguientes métodos devuelven el número de errores hasta el momento (`YYNberr`) y si nos encontramos o no en fase de recuperación de errores (`YYRecovering`):

```

sub YYNberr {
    my($self)=shift;

    ${$self{NBERR}};
}

sub YYRecovering {
    my($self)=shift;

    ${$self{ERRST}} != 0;
}

```

7.27. Descripción Eyapp del Lenguaje SimpleC

En este capítulo usaremos `Parse::Eyapp` para desarrollar un compilador para el siguiente lenguaje, al que denominaremos `Simple C`:

```

program: definition+

definition: funcDef | basictype funcDef | declaration

basictype: INT | CHAR

funcDef: ID '(' params ')' block

params: ( basictype ID arraySpec )<* ','>

block: '{' declaration* statement* '}'

```



```

declaration: basictype declList ';'

declList: (ID arraySpec) <+ ','>

arraySpec: ( '[' INUM ''])*

statement:
    expression ';'
  | ';'
  | BREAK ';'
  | CONTINUE ';'
  | RETURN ';'
  | RETURN expression ';'
  | block
  | ifPrefix statement %prec '+'
  | ifPrefix statement 'ELSE' statement
  | loopPrefix statement

ifPrefix: IF '(' expression ')'

loopPrefix: WHILE '(' expression ')'

expression: binary <+ ','>

Variable: ID ( '[' binary ''] ) *

Primary:
    INUM
  | CHARCONSTANT
  | Variable
  | '(' expression ')'
  | function_call

function_call: ID '(' binary <* ','> ')'

Unary: '++' Variable | '--' Variable | Primary

binary:
    Unary
  | binary '+' binary
  | binary '-' binary
  | binary '*' binary
  | binary '/' binary
  | binary '%' binary
  | binary '<' binary
  | binary '>' binary
  | binary '>=' binary
  | binary '<=' binary
  | binary '==' binary
  | binary '!=' binary
  | binary '&' binary
  | binary '**' binary
  | binary '|' binary

```

```
| Variable '=' binary
| Variable '+=' binary
| Variable '-=' binary
| Variable '*=' binary
| Variable '/=' binary
| Variable '%=' binary
| etc. etc.
```

7.28. Diseño de Analizadores con Parse::Eyapp

A la hora de construir un analizador sintáctico tenga en cuenta las siguientes normas de buena programación:

1. Comience trabajando en el cuerpo de la gramática.
2. Olvídense al principio del analizador léxico. Su primer objetivo es tener una gramática limpia de conflictos y que reconozca el lenguaje dado.
3. Sustituya las repeticiones BNF por listas usando los operadores **eyapp** `+`, `*` y sus variantes con separadores. Si una variable describe una lista de **cosas** pongale un adjetivo adecuado como **cosaslist**. Ponga nombres significativos a las variables y terminales. No los llame **d1**, **d2**, etc.
4. Si tiene un elemento opcional en la BNF, por ejemplo, en la regla:
`functiondefinition` \rightarrow `[basictype] functionheader functionbody`
use el operador `?`.
5. Cada par de reglas que introduzca vuelva a recompilar con **eyapp** la gramática para ver si se introducido ambigüedad. Cuando estoy editando la gramática suelo escribir a menudo la orden
`:!eyapp %`
para recompilar:

[illegible]

Esto llama a **eyapp** con el fichero bajo edición. Si hay errores o conflictos (esto es, hemos introducido ambigüedad) los detectaremos enseguida. Procure detectar la aparición de un conflicto lo antes posible. Observe el sangrado del ejemplo. Es el que le recomiendo.

6. Cuando esté en el proceso de construcción de la gramática y aún le queden por rellenar variables sintácticas, declárelas como terminales mediante **%token**. De esta manera evitará las quejas de **eyapp**.

7. Resolución de Ambigüedades y Conflictos

Las operaciones de asignación tienen la prioridad mas baja, seguidas de las lógicas, los test de igualdad, los de comparación, a continuación las aditivas, multiplicativas y por último las operaciones de tipo **unary** y **primary**. Expresa la asociatividad natural y la prioridad especificada usando los mecanismos que **eyapp** provee al efecto: **%left**, **%right**, **%nonassoc** y **%prec**.

8. La gramática de SimpleC es ambigua, ya que para una sentencia como

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$$

existen dos árboles posibles: uno que asocia el “else” con el primer “if” y otra que lo asocia con el segundo. Los dos árboles corresponden a las dos posibles parentizaciones:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1 \text{ else } S_2)$$

Esta es la regla de prioridad usada en la mayor parte de los lenguajes: un “else” casa con el “if” mas cercano. La otra posible parentización es:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1) \text{ else } S_2$$

*La conducta por defecto de **eyapp** es parentizar a derechas.* El generador **eyapp** nos informará del conflicto pero si no se le indica como resolverlo parentizará a derechas. Resuelva este conflicto.

9. *¿Que clase de árbol debe producir el analizador?* La respuesta es que sea lo mas abstracto posible. Debe

- Contener toda la información necesaria para el manejo eficiente de las fases subsiguientes: Análisis de ámbito, Comprobación de tipos, Optimización independiente de la máquina, etc.
- Ser uniforme
- Legible (human-friendly)
- No contener nodos que no portan información.

El siguiente ejemplo muestra una versión aceptable de árbol abstracto. Cuando se le proporciona el programa de entrada:

```
nereida:~/doc/casiano/PLB00K/PLB00K/code> cat -n prueba5.c
1  int f(int a)
2  {
3      if (a>0)
4          a = f(a-1);
5  }
```

El siguiente árbol ha sido producido por un analizador usando la directiva **%tree** y añadiendo las correspondientes acciones de **bypass**. Puede considerarse un ejemplo aceptable de AST:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> eyapp Simple2 ;\
                                         usesimple2.pl prueba5.c
```

```
PROGRAM(
  TYPEDFUNC(
    INT(TERMINAL[INT:1]),
    FUNCTION(
      TERMINAL[f:1],
      PARAMS(
        PARAM(
          INT(TERMINAL[INT:1]),
          TERMINAL[a:1],
          ARRAYSPEC
        )
      ),
      BLOCK(
        DECLARATIONS,
        STATEMENTS(
          IF(
            GT(
              VAR(TERMINAL[a:3]),
              INUM(TERMINAL[0:3])
            ),
            ASSIGN(
              VAR(TERMINAL[a:4]),
              FUNCTIONCALL(
                TERMINAL[f:4],
                ARGLIST(
                  MINUS(
                    VAR(TERMINAL[a:4]),
                    INUM(TERMINAL[1:4])
                  )
                ) # ARGLIST
              ) # FUNCTIONCALL
            ) # ASSIGN
          ) # IF
        ) # STATEMENTS
      ) # BLOCK
    ) # FUNCTION
  ) # TYPEDFUNC
) # PROGRAM
```

Es deseable darle una estructura uniforme al árbol. Por ejemplo, como consecuencia de que la gramática admite funciones con declaración implícita del tipo retornado cuando este es entero

```
1 definition:
2   funcDef { $_[1]->type("INTFUNC"); $_[1] }
3   | %name TYPEDFUNC
4   basictype funcDef
5   | declaration { $_[1] }
6   ;
```

se producen dos tipos de árboles. Es conveniente convertir las definiciones de función con declaración implícita en el mismo árbol que se obtiene con declaración explícita.

7.29. Práctica: Construcción del AST para el Lenguaje Simple C

Utilice Parse-Eyapp. para construir un árbol de análisis sintáctico abstracto para la gramática descrita en la sección 12.2. Su analizador deberá seguir los consejos explícitados en la sección 8.17.

Analizador Léxico Además del tipo de terminal y su valor el analizador léxico deberá devolver el número de línea. El analizador léxico deberá aceptar comentarios C. En la gramática, el terminal **CHARACTER** se refiere a caracteres entre comillas simples (por ejemplo 'a'). El terminal **STRING** se refiere a caracteres entre comillas dobles (por ejemplo "hola").

Se aconseja que las palabras reservadas del lenguaje no se traten con expresiones regulares específicas sino que se capturen en el patrón de identificador `[a-z_]\w+`. Se mantiene para ello un hash con las palabras reservadas que es inicializado al comienzo del programa. Cuando el analizador léxico encuentra un identificador mira en primer lugar en dicho hash para ver si es una palabra reservada y, si lo es, devuelve el terminal correspondiente. En caso contrario se trata de un identificador.

7.30. El Generador de Analizadores byacc

Existe una version del yacc de Berkeley que permite producir código para Perl:

```
> byacc -V
byacc: Berkeley yacc version 1.8.2 (C or perl)
```

Se trata por tanto de un generador de analizadores LALR. Es bastante compatible con AT&T yacc. Puedes encontrar una versión en formato tar.gz en nuestro servidor <http://nereida.deioc.ull.es/~pl/pyacc-pack.tgz> o también desde <http://www.perl.com/CPAN/src/misc/>.

El formato de llamada es:

```
byacc [ -CPcdlrtv ] [ -b file_prefix ] [ -p symbol_prefix ] filename
```

Las opciones C o c permiten generar código C. Usando -P se genera código Perl. Las opciones d y v funcionan como es usual en yacc. Con t se incorpora código para la depuración de la gramática. Si se especifica l el código del usuario no es insertado. La opción r permite generar ficheros separados para el código y las tablas. No la use con Perl.

Fichero conteniendo la gramática:

```
%{
%}

%token INT EOL
%token LEFT_PAR RIGHT_PAR
%left PLUS MINUS
%left MULT DIV

%%
start: |
start input
;

input: expr EOL { print $1 . "\n"; }
| EOL
;

expr: INT { $p->mydebug("INT -> Expr!"); $$ = $1; }
| expr PLUS expr { $p->mydebug("PLUS -> Expr!"); $$ = $1 + $3; }
| expr MINUS expr { $p->mydebug("MINUS -> Expr!"); $$ = $1 - $3; }
```

```

| expr MULT expr { $p->mydebug("MULT -> Expr!"); $$ = $1 * $3; }
| expr DIV expr { $p->mydebug("DIV -> Expr!"); $$ = $1 / $3; }
| LEFT_PAR expr RIGHT_PAR { $p->mydebug("PARENS -> Expr!"); $$ = $2; }
;
%%

sub yyerror {
    my ($msg, $s) = @_;
    my ($package, $filename, $line) = caller;

    die "$msg at <DATA> \n$package\n$filename\n$line\n";
}

sub mydebug {
    my $p = shift;
    my $msg = shift;
    if ($p->{'yydebug'})
    {
        print "$msg\n";
    }
}

```

La compilación con `byacc` del fichero `calc.y` conteniendo la descripción de la gramática produce el módulo Perl conteniendo el analizador.

```

> ls -l
total 12
-rw-r----- 1 pl      casiano      47 Dec 29 2002 Makefile
-rw-r----- 1 pl      casiano      823 Dec 29 2002 calc.y
-rwxr-x--x  1 pl      casiano      627 Nov 10 15:37 tokenizer.pl
> cat Makefile
MyParser.pm: calc.y
        byacc -d -P MyParser $<

> make
byacc -d -P MyParser calc.y
> ls -ltr
total 28
-rw-r----- 1 pl      casiano      823 Dec 29 2002 calc.y
-rw-r----- 1 pl      casiano      47 Dec 29 2002 Makefile
-rwxr-x--x  1 pl      casiano      627 Nov 10 15:37 tokenizer.pl
-rw-rw----  1 pl      users        95 Nov 16 12:49 y.tab.ph
-rw-rw----  1 pl      users      9790 Nov 16 12:49 MyParser.pm

```

Observe que la opción `-P` es la que permite producir código Perl. Anteriormente se usaba la opción `-p`. Esto se hizo para mantener la compatibilidad con otras versiones de `yacc` en las que la opción `-p` se usa para cambiar el prefijo por defecto (`yy`). Ese es el significado actual de la opción `-p` en `perl-byacc`.

El fichero `y.tab.ph` generado contiene las definiciones de los *tokens*:

```

cat y.tab.ph
$INT=257;
$EOL=258;
$LEFT_PAR=259;
$RIGHT_PAR=260;

```

```
$PLUS=261;
$MINUS=262;
$MULT=263;
$DIV=264;
```

El programa `tokenizer.pl` contiene la llamada al analizador y la definición del analizador léxico:

```
> cat tokenizer.pl
#!/usr/local/bin/perl5.8.0

require 5.004;
use strict;
use Parse::YYLex;
use MyParser;

print STDERR "Version $Parse::ALex::VERSION\n";

my (@tokens) = ((LEFT_PAR => '\(',
                 RIGHT_PAR => '\)',
                 MINUS => '-',
                 PLUS => '+',
                 MULT => '*',
                 DIV => '/',
                 INT => '[1-9][0-9]*',
                 EOL => '\n',
                 ERROR => '.*'),
               sub { die "!can't analyze: \"$_[1]\"\n!"; });

my $lexer = Parse::YYLex->new(@tokens);

sub yyerror
{
    die "There was an error:" . join("\n", @_). "\n";
}

my $debug = 0;
my $parser = new MyParser($lexer->getyylex(), \&MyParser::yyerror, $debug);
$lexer->from(\*STDIN);
$parser->yyparse(\*STDIN);
```

El módulo `Parse::YYLex` contiene una versión de `Parse::Lex` que ha sido adaptada para funcionar con `byacc`. Todas las versiones de `yacc` esperan que el analizador léxico devuelva un *token* numérico, mientras que `Parse::Lex` devuelve un objeto de la clase *token*. Veamos un ejemplo de ejecución:

```
> tokenizer.pl
Version 2.15
yydebug: state 0, reducing by rule 1 (start :)
yydebug: after reduction, shifting from state 0 to state 1
3*(5-9)
yydebug: state 1, reading 257 (INT)
yydebug: state 1, shifting to state 2
yydebug: state 2, reducing by rule 5 (expr : INT)
INT -> Expr!
yydebug: after reduction, shifting from state 1 to state 6
yydebug: state 6, reading 263 (MULT)
```

```

yydebug: state 6, shifting to state 11
yydebug: state 11, reading 259 (LEFT_PAR)
yydebug: state 11, shifting to state 4
yydebug: state 4, reading 257 (INT)
yydebug: state 4, shifting to state 2
yydebug: state 2, reducing by rule 5 (expr : INT)
INT -> Expr!
yydebug: after reduction, shifting from state 4 to state 7
yydebug: state 7, reading 262 (MINUS)
yydebug: state 7, shifting to state 10
yydebug: state 10, reading 257 (INT)
yydebug: state 10, shifting to state 2
yydebug: state 2, reducing by rule 5 (expr : INT)
INT -> Expr!
yydebug: after reduction, shifting from state 10 to state 15
yydebug: state 15, reading 260 (RIGHT_PAR)
yydebug: state 15, reducing by rule 7 (expr : expr MINUS expr)
MINUS -> Expr!
yydebug: after reduction, shifting from state 4 to state 7
yydebug: state 7, shifting to state 13
yydebug: state 13, reducing by rule 10 (expr : LEFT_PAR expr RIGHT_PAR)
PARENS -> Expr!
yydebug: after reduction, shifting from state 11 to state 16
yydebug: state 16, reducing by rule 8 (expr : expr MULT expr)
MULT -> Expr!
yydebug: after reduction, shifting from state 1 to state 6
yydebug: state 6, reading 258 (EOL)
yydebug: state 6, shifting to state 8
yydebug: state 8, reducing by rule 3 (input : expr EOL)
-12
yydebug: after reduction, shifting from state 1 to state 5
yydebug: state 5, reducing by rule 2 (start : start input)
yydebug: after reduction, shifting from state 0 to state 1
yydebug: state 1, reading 0 (end-of-file)

```


Capítulo 8

Análisis Sintáctico con Parse::Eyapp

8.1. Conceptos Básicos para el Análisis Sintáctico

Suponemos que el lector de esta sección ha realizado con éxito un curso en teoría de autómatas y lenguajes formales. Las siguientes definiciones repasan los conceptos mas importantes.

Definición 8.1.1. Dado un conjunto A , se define A^* el cierre de Kleene de A como: $A^* = \bigcup_{n=0}^{\infty} A^n$. Se admite que $A^0 = \{\epsilon\}$, donde ϵ denota la palabra vacía, esto es la palabra que tiene longitud cero, formada por cero símbolos del conjunto base A .

Definición 8.1.2. Una gramática G es una cuaterna $G = (\Sigma, V, P, S)$. Σ es el conjunto de terminales. V es un conjunto (disjunto de Σ) que se denomina conjunto de variables sintácticas o categorías gramaticales, P es un conjunto de pares de $V \times (V \cup \Sigma)^*$. En vez de escribir un par usando la notación $(A, \alpha) \in P$ se escribe $A \rightarrow \alpha$. Un elemento de P se denomina producción. Por último, S es un símbolo del conjunto V que se denomina símbolo de arranque.

Definición 8.1.3. Dada una gramática $G = (\Sigma, V, P, S)$ y $\mu = \alpha A \beta \in (V \cup \Sigma)^*$ una frase formada por variables y terminales y $A \rightarrow \gamma$ una producción de P , decimos que μ deriva en un paso en $\alpha \gamma \beta$. Esto es, derivar una cadena $\alpha A \beta$ es sustituir una variable sintáctica A de V por la parte derecha γ de una de sus reglas de producción. Se dice que μ deriva en n pasos en δ si deriva en $n - 1$ pasos en una cadena η la cual deriva en un paso en δ . Se escribe entonces que $\mu \xRightarrow{n} \delta$. Una cadena deriva en 0 pasos en si misma. Se escribe entonces que $\mu \xRightarrow{*} \delta$.

Definición 8.1.4. Dada una gramática $G = (\Sigma, V, P, S)$ se denota por $L(G)$ o lenguaje generado por G al lenguaje:

$$L(G) = \{x \in \Sigma^* : S \xRightarrow{*} x\}$$

Esto es, el lenguaje generado por la gramática G esta formado por las cadenas de terminales que pueden ser derivados desde el símbolo de arranque.

Definición 8.1.5. Una derivación que comienza en el símbolo de arranque y termina en una secuencia formada por sólo terminales de Σ se dice completa.

Una derivación $\mu \xRightarrow{*} \delta$ en la cual en cada paso $\alpha A x$ la regla de producción aplicada $A \rightarrow \gamma$ se aplica en la variable sintáctica mas a la derecha se dice una derivación a derechas

Una derivación $\mu \xRightarrow{*} \delta$ en la cual en cada paso $x A \alpha$ la regla de producción aplicada $A \rightarrow \gamma$ se aplica en la variable sintáctica mas a la izquierda se dice una derivación a izquierdas

Definición 8.1.6. Observe que una derivación puede ser representada como un árbol cuyos nodos están etiquetados en $V \cup \Sigma$. La aplicación de la regla de producción $A \rightarrow \gamma$ se traduce en asignar como hijos del nodo etiquetado con A a los nodos etiquetados con los símbolos $X_1 \dots X_n$ que constituyen la frase $\gamma = X_1 \dots X_n$. Este árbol se llama árbol sintáctico concreto asociado con la derivación.

Definición 8.1.7. Observe que, dada una frase $x \in L(G)$ una derivación desde el símbolo de arranque da lugar a un árbol. Ese árbol tiene como raíz el símbolo de arranque y como hojas los terminales $x_1 \dots x_n$ que forman x . Dicho árbol se denomina árbol de análisis sintáctico concreto de x . Una derivación determina una forma de recorrido del árbol de análisis sintáctico concreto.

Definición 8.1.8. Una gramática G se dice ambigua si existe alguna frase $x \in L(G)$ con al menos dos árboles sintácticos. Es claro que esta definición es equivalente a afirmar que existe alguna frase $x \in L(G)$ para la cual existen dos derivaciones a izquierda (derecha) distintas.

Definición 8.1.9. Un esquema de traducción es una gramática independiente del contexto en la cual se han insertado fragmentos de código en las partes derechas de sus reglas de producción.

$$A \rightarrow \alpha\{action\}$$

Los fragmentos de código así insertados se denominan acciones semánticas.

En un esquema de traducción los nodos del árbol sintáctico tienen asociados atributos. Si pensamos que cada nodo del árbol es un objeto, entonces los atributos del nodo son los atributos del objeto. Las reglas semánticas determinan la forma en la que son evaluados los atributos.

Los fragmentos de código de un esquema de traducción calculan y modifican los atributos asociados con los nodos del árbol sintáctico. El orden en que se evalúan los fragmentos es el de un recorrido primero-profundo del árbol de análisis sintáctico. Esto significa que si en la regla $A \rightarrow \alpha\beta$ insertamos un fragmento de código:

$$A \rightarrow \alpha\{action\}\beta$$

La acción $\{action\}$ se ejecutará después de todas las acciones asociadas con el recorrido del subárbol de α y antes que todas las acciones asociadas con el recorrido del subárbol β .

Obsérvese que para poder aplicar un esquema de traducción hay que - al menos conceptualmente - construir el árbol sintáctico y después aplicar las acciones empotradas en las reglas en el orden de recorrido primero-profundo. Por supuesto, si la gramática es ambigua una frase podría tener dos árboles y la ejecución de las acciones para ellos podría dar lugar a diferentes resultados. Si se quiere evitar la multiplicidad de resultados (interpretaciones semánticas) es necesario precisar de que árbol sintáctico concreto se está hablando.

Definición 8.1.10. Un atributo tal que su valor en un nodo puede ser computado en términos de los atributos de los hijos del nodo se dice que es un atributo sintetizado.

El siguiente esquema de traducción recibe como entrada una expresión en infijo y produce como salida su traducción a postfijo para expresiones aritmeticas con sólo restas de números:

$$\begin{array}{ll} expr \rightarrow expr_1 - NUM & \{ \$expr\{TRA\} = \$expr[1]\{TRA\}." " . \$NUM\{VAL\}." - " \} \\ expr \rightarrow NUM & \{ \$expr\{TRA\} = \$NUM\{VAL\} \} \end{array}$$

Que para una entrada $2 - 3 - 7$ daría lugar a la siguiente evaluación:

```
e '2 3 - 7 -'
'-- --- e '2 3 -'
|   '----- e '2'
|       |   '-- N 2
|       |   '-- '-'
|       |   '-- N 3
|       |   '-- '-'
|       |   '-- N 7
```

Definición 8.1.11. Un atributo heredado es aquel cuyo valor se computa a partir de los valores de sus hermanos y de su padre.

Ejemplo 8.1.1. Un ejemplo de atributo heredado es el tipo de las variables en las declaraciones:

```
decl → type { $list{T} = $type{T} } list
type → INT { $type{T} = $int }
type → STRING { $type{T} = $string }
list → ID , { $ID{T} = $list{T}; $list_1{T} = $list{T} } list_1
list → ID { $ID{T} = $list{T} }
```

8.2. Parse::Eyapp: Un Generador de Analizadores Sintácticos

El generador de analizadores sintácticos `Parse::Eyapp` es un analizador LALR inspirado en `yacc`. `Parse::Eyapp` es una extensión de `Parse::Yapp` escrita por Casiano Rodríguez-Leon. Puede descargarlo desde CPAN e instalarlo en su máquina siguiendo el procedimiento habitual.

El generador de analizadores sintácticos `Parse::Eyapp` que estudiaremos en las siguientes secciones funciona de manera similar a un esquema de traducción. Las reglas de producción de la gramática son aumentadas con reglas semánticas. Los símbolos que aparecen en la regla de producción tienen *atributos* asociados y las reglas semánticas dicen como deben ser computados dichos atributos. Consideremos, por ejemplo, el siguiente fragmento de programa `eyapp`:

```
exp:  exp '+' exp          { $_[1] + $_[3] }
```

que dice que asociado con el símbolo `exp` de la regla de producción $exp \rightarrow exp' + exp$ tenemos el atributo valor numérico y que para computar el atributo valor de la variable sintáctica `exp` en la parte izquierda tenemos que sumar los atributos asociados con los símbolos primero y tercero de la parte derecha.

Por defecto `Parse::Eyapp` no provee un esquema de traducción ya que - aunque el orden de ejecución de las acciones es de abajo-arriba y de izquierda a derecha como en un esquema - no es posible acceder a atributos de nodos que aún no han sido visitados.

Para ilustrar el uso de `Parse::Eyapp` veamos un ejemplo en el que se implanta una gramática cuyas frases son secuencias (separadas por retornos de carro) de expresiones aritméticas.

Los contenidos del programa `eyapp` los hemos guardado en un fichero denominado `CalcSyntax.eyp`

Partes de un Programa Eyapp

Un programa `eyapp` consta de tres partes:

- la cabeza,
- el cuerpo
- y la cola.

Cada una de las partes va separada de las otras por el símbolo `%` en una línea aparte. Así, el `%` de la línea 8 separa la cabeza del cuerpo y el de la línea 31 el cuerpo de la cola.

- En la cabecera se colocan el código de inicialización, las declaraciones de terminales, las reglas de precedencia, etc.
- El cuerpo contiene las reglas de la gramática y las acciones asociadas.
- Por último, la cola de un program `eyapp` contiene las rutinas de soporte al código que aparece en las acciones así como, posiblemente, rutinas para el análisis léxico y el tratamiento de errores.

```
pl@nereida:~/LEyapp/examples$ cat -n CalcSyntax.eyy
```

```

1  # CalcSyntax.eyy
2  %right  '='
3  %left   '-' '+'
4  %left   '*' '/'
5  %left   NEG
6  %right  '^'
7
8  %%
9
10 input:  line * { print "input -> line *\n" }
11 ;
12
13 line:
14     '\n'          { print "line -> \n\n" }
15     | exp '\n'    { print "line -> exp \n\n" }
16 ;
17
18 exp:
19     NUM            { print "exp -> NUM ($_[1])\n"; }
20     | VAR          { print "exp -> VAR ($_[1])\n"; }
21     | VAR '=' exp  { print "exp -> VAR '=' exp\n"; }
22     | exp '+' exp  { print "exp -> exp '+' exp\n"; }
23     . .....
29 ;
30
31 %%
32
33 sub _Error {
34     .....
41 }
42
43 .....
44
45 sub _Lexer {
46     my($parser)=shift;
47     .....
58 }
59
60 sub Run {
61     my($self)=shift;
62
63     $input = shift;
64     return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
65 }

```

Las Reglas

Todas las partes derechas de las reglas de producción de una misma variable sintáctica se escriben juntas separadas mediante la barra vertical |.

```

10 input:  line * { print "input -> line *\n" }
11 ;
12
13 line:

```

```

14     '\n'          { print "line -> \\n\\n" }
15     | exp '\n'    { print "line -> exp \\n\\n"}
16 ;

```

En este ejemplo hemos simplificado las acciones semánticas reduciéndolas a mostrar la regla de producción encontrada.

Reglas de Producción Vacías

Un asterisco (como en la línea 10) indica repetición cero o mas veces de la expresión a la que se aplica. De hecho la línea 10 es casi equivalente a:

```

input: temp
;
temp:
    /* vacio */
    | temp line
;

```

Observe como en el código anterior hemos codificado la regla de producción $temp \rightarrow \epsilon$ como:

```

temp:
    /* vacio */

```

Es buena costumbre de programación cuando se tiene una regla de producción que produce vacío ponerla la primera del grupo y añadir un comentario como este. Dado que vacío se representa en **Eyapp** mediante la cadena vacía es fácil que pase desapercibida. Es por ello que se recomienda que *una regla vacía sea siempre la primera y que este comentada* como en el ejemplo.

Tratamiento de las Ambigüedades

Hay numerosas ambigüedades en esta gramática. Observe las reglas para los diferentes tipos de expresiones:

```

18 exp:
19     NUM                { print "exp -> NUM ($_[1])\n"; }
20     | VAR              { print "exp -> VAR ($_[1])\n"; }
21     | VAR '=' exp      { print "exp -> VAR '=' exp\n"; }
22     | exp '+' exp      { print "exp -> exp '+' exp\n"; }
23     | exp '-' exp      { print "exp -> exp '-' exp\n"; }
24     | exp '*' exp      { print "exp -> exp '*' exp\n"; }
25     | exp '/' exp      { print "exp -> exp '/' exp\n"; }
26     | '-' exp %prec NEG { print "exp -> '-' exp\n"; }
27     | exp '^' exp      { print "exp -> exp '^' exp\n"; }
28     | '(' exp ')'      { print "exp -> '(' exp ')' \n"; }
29 ;

```

Surgen preguntas como:

- ¿Como debo interpretar la expresión $4 - 5 - 2$? ¿Como $(4 - 5) - 2$? ¿o bien $4 - (5 - 2)$? La respuesta la da la asignación de asociatividad a los operadores que hicimos en la cabecera:

```

1 # CalcSyntax.eypp
2 %right  '='
3 %left  '-' '+'
4 %left  '*' '/'
5 %left  NEG
6 %right '^'
7
8 %%

```

La Asociatividad de Terminales y la Ambigüedad

Las declaraciones `%left` y `%right` expresan la asociatividad y precedencia de los terminales, permitiendo decidir que árbol construir en caso de ambigüedad.

Los terminales declarados en líneas posteriores tienen mas prioridad que los declarados en las líneas anteriores.

Por defecto, una regla de producción tiene la prioridad del último terminal que aparece en su parte derecha.

Al declarar como asociativo a izquierdas al terminal `-` hemos resuelto la ambigüedad en $4 - 5 - 2$. Lo que estamos haciendo es indicarle al analizador que a la hora de elegir entre los árboles abstractos $-(4, 5), 2$ y $-(4, -(5, 2))$ elija siempre el árbol que se hunde a izquierdas.

- ¿Como debo interpretar la expresión $4 - 5 * 2$? ¿Como $(4 - 5) * 2$? ¿o bien $4 - (5 * 2)$? Al declarar que `*` tiene mayor prioridad que `-` estamos resolviendo esta otra fuente de ambigüedad. Esto es así pues `*` fué declarado en la línea 11 y `-` en la 10. Por tanto el árbol será $-(4, *(5, 2))$.

La declaración de `^` como asociativo a derechas y con un nivel de prioridad alto resuelve las ambigüedades relacionadas con este operador:

```
18  exp:
..      .....
28      | '(' exp ')',          { print "exp -> '(' exp ')'\n"; }
```

Ejercicio 8.2.1. ¿Como se esta interpretando la expresión -2^2 ? ¿Cómo $(-2)^2$? ¿o bien $-(2^2)$?

Modificación de la Prioridad Implícita de una Regla

Una regla de producción puede ir seguida de una directiva `%prec` la cual le da una prioridad explícita. Esto puede ser de gran ayuda en ciertos casos de ambigüedad.

```
26      | '-' exp %prec NEG { print "exp -> '-' exp\n"; }
```

¿Cual es la ambigüedad que surge con esta regla? La ambigüedad de esta regla esta relacionada con el doble significado del menos como operador unario y binario: hay frases como $-y-z$ que tiene dos posibles interpretaciones: Podemos verla como $(-y)-z$ o bien como $-(y-z)$. Hay dos árboles posibles. El analizador, cuando este analizando la entrada $-y-z$ y vea el segundo `-` (después de haber leído $-y$) deberá escoger uno de los dos árboles. ¿Cuál?. El conflicto puede verse como una “lucha” entre la regla `exp: '-' exp` la cual interpreta la frase como $(-y)-z$ y la segunda aparición del terminal `-` el cual “quiere entrar” para que gane la regla `exp: exp '-' exp` y dar lugar a la interpretación $-(y-z)$. En este caso, las dos reglas $E \rightarrow -E$ y $E \rightarrow E - E$ tienen, en principio la prioridad del terminal `-`, el cual fué declarado en la zona de cabecera:

```
1 # CalcSyntax.eypp
2 %right  '='
3 %left   '-' '+'
4 %left   '*' '/'
5 %left   NEG
6 %right  '^'
7
8 %%
```

La prioridad expresada explícitamente para la regla por la declaración `%prec NEG` de la línea 41 hace que la regla tenga la prioridad del terminal `NEG` y por tanto mas prioridad que el terminal `-`. Esto hará que `eyapp` finalmente opte por la regla `exp: '-' exp` dando lugar a la interpretación $(-y)-z$.

La Cola

Después de la parte de la gramática, y separada de la anterior por el símbolo `%%`, sigue la parte en la que se suelen poner las rutinas de apoyo. Hay al menos dos rutinas de apoyo que el analizador sintáctico requiere le sean pasados como argumentos: la de manejo de errores y la de análisis léxico.

```
31  %%
32
33  sub _Error {
34      exists $_[0]->YYData->{ERRMSG}
35      and do {
36          print $_[0]->YYData->{ERRMSG};
37          delete $_[0]->YYData->{ERRMSG};
38          return;
39      };
40      print "Syntax error.\n";
41  }
42
43  my $input;
44
45  sub _Lexer {
46      my($parser)=shift;
47
48      # topicalize $input
49      for ($input) {
50          s/^[ \t]+//;      # skip whites
51          return('','undef) unless $input;
52
53          return('NUM',$1) if s/^[0-9]+(?:\.[0-9]+)?/{}/;
54          return('VAR',$1) if s/^[A-Za-z][A-Za-z0-9_]*/;/;
55          return($1,$1)    if s/^(.)/s/;
56      }
57  }
58
59  sub Run {
60      my($self)=shift;
61
62      $input = shift;
63      return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
64  }
```

El Método Run

El método `Run` ilustra como se hace la llamada al método de análisis sintáctico generado, utilizando la técnica de llamada con argumentos con nombre y pasándole las referencias a las dos subrutinas (en Perl, es un convenio que si el nombre de una subrutina comienza por un guión bajo es que el autor la considera privada):

```
78  sub Run {
79      my($self)=shift;
80
81      $input = shift;
82      return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
83  }
```

El Atributo YYData y el Método _Error

El método `YYData` provee acceso a un hash que contiene los datos que están siendo analizados.

La subrutina de manejo de errores `_Error` imprime el mensaje de error proveído por el usuario, el cual, si existe, fué guardado en `$_[0]->YYData->{ERRMSG}`.

```
51 sub _Error {
52     exists $_[0]->YYData->{ERRMSG}
53     and do {
54         print $_[0]->YYData->{ERRMSG};
55         delete $_[0]->YYData->{ERRMSG};
56         return;
57     };
58     print "Syntax error.\n";
```

El Método `_Lexer`

A continuación sigue el método que implanta el análisis léxico `_Lexer`.

```
45 sub _Lexer {
46     my($parser)=shift;
47
48     # topicalize $input
49     for ($input) {
50         s/^[ \t]+//;      # skip whites
51         return('','undef') unless $input;
52
53         return('NUM',$1) if s/^([0-9]+(?:\.[0-9]+)?)\b/;
54         return('VAR',$1) if s/^[A-Za-z][A-Za-z0-9_]*\b/;
55         return($1,$1)    if s/^\b/;
56     }
```

El bucle `for` constituye en este caso una "frase hecha": el efecto es hacer que `$_` sea un alias de `$input`. Se simplifica la escritura (obsérvese que no es necesario explicitar el operador de binding en las líneas 50-55, no teniendo que escribir `$input =~ s/^([0-9]+(?:\.[0-9]+)?)\b/`) y que la computación será mas eficiente al acceder a través de `$_` en vez de `$input`.

El bucle `for ($input)` se ejecutará mientras la cadena en `$input` no sea vacía, lo que ocurrirá cuando todos los terminales hayan sido consumidos. sin embargo es un "falso for": no hay iteración. El interior del bucle es ejecutado una sola vez en cada llamada.

Eliminamos los blancos iniciales (lo que en inglés se conoce por *trimming*) y a continuación vamos detectando los números, identificadores y los símbolos individuales.

En primer lugar se comprueba la existencia de datos. Si no es el caso, estamos ante el final de la entrada. Cuando el analizador léxico alcanza el final de la entrada debe devolver la pareja `('','undef')`.

Ejercicio 8.2.2. 1. ¿Quién es la variable índice en la línea 49?

2. ¿Sobre quién ocurre el binding en las líneas 50-54?

3. ¿Cual es la razón por la que `$input` se ve modificado aún cuando no aparece como variable para el binding en las líneas 50-54?

4. Dada la entrada `'4 * 3 '` con blancos al final: como termina el analizador léxico. ¿Funciona correctamente en ese caso?

Compilación con `eyapp`

Construimos el módulo `CalcSyntax.pm` a partir del fichero `CalcSyntax.eyp` especificando la gramática, usando ejecutable `eyapp`:


```
pl@nereida:~/LEyapp/examples$ eyapp -m CalcSyntax CalcSyntax.eyp
pl@nereida:~/LEyapp/examples$ ls -ltr | tail -3
-rw-r--r-- 1 pl users 1545 2007-10-24 09:03 CalcSyntax.eyp
-rwxr-xr-x 1 pl users 329 2007-10-24 09:05 usecalcsyntax.pl
-rw-r--r-- 1 pl users 7848 2007-10-24 09:36 CalcSyntax.pm
```

Esta compilación genera el fichero `CalcSyntax.pm` conteniendo el analizador.

El script `eyapp` es un *frontend* al módulo `Parse::Eyapp`. Admite diversas formas de uso:

- `eyapp [options] grammar[.eyp]`

Los sufijos `.eyp` io `.yp` son opcionales.

- `eyapp -V`

Nos muestra la versión:

```
pl@nereida:~/LEyapp/examples$ eyapp -V
This is Parse::Eyapp version 1.081.
```

- `eyapp -h`

Nos muestra la ayuda:

```
pl@nereida:~/LEyapp/examples$ eyapp -h
```

```
Usage: eyapp [options] grammar[.yp]
or eyapp -V
or eyapp -h
```

```
-m module    Give your parser module the name <module>
              default is <grammar>
-v           Create a file <grammar>.output describing your parser
-s           Create a standalone module in which the driver is included
-n           Disable source file line numbering embedded in your parser
-o outfile   Create the file <outfile> for your parser module
              Default is <grammar>.pm or, if -m A::Module::Name is
              specified, Name.pm
-t filename  Uses the file <filename> as a template for creating the parser
              module file. Default is to use internal template defined
              in Parse::Eyapp::Output
-b shebang   Adds '#!<shebang>' as the very first line of the output file

grammar      The grammar file. If no suffix is given, and the file
              does not exists, .yp is added

-V           Display current version of Parse::Eyapp and gracefully exits
-h           Display this help screen
```

La opción `-o outfile` da el nombre del fichero de salida. Por defecto toma el nombre de la gramática, seguido del sufijo `.pm`. sin embargo, si hemos especificado la opción `-m A::Module::Name` el valor por defecto será `Name.pm`.

La Jerarquía de un Módulo y La Opción *-m module*

La opción *-m module* da el nombre al paquete o espacio de nombres o clase encapsulando el analizador. Por defecto toma el nombre de la gramática. En el ejemplo anterior podría haberse omitido. Sin embargo es necesaria cuando se está desarrollando un módulo con un nombre complejo. Construyamos una distribución con `h2xs`:

```
$ h2xs -XA -n Calc::Syntax
Writing Calc-Syntax/lib/Calc/Syntax.pm
Writing Calc-Syntax/Makefile.PL
Writing Calc-Syntax/README
Writing Calc-Syntax/t/Calc-Syntax.t
Writing Calc-Syntax/Changes
Writing Calc-Syntax/MANIFEST
```

Ahora añadimos el fichero `.eypp` en el directorio de la librería y producimos el módulo `Syntax.pm` al compilarlo. Para darle al paquete el nombre `Calc::Syntax` usamos la opción `-m`:

```
$ cd Calc-Syntax/lib/Calc/
$ cp ~/LEyapp/examples/CalcSyntax.eypp .
$ eyapp -m Calc::Syntax CalcSyntax.eypp
$ head -12 Syntax.pm | cat -n
  1 #####
  2 #
  3 #   This file was generated using Parse::Eyapp version 1.081.
  4 #
  5 # (c) Parse::Yapp Copyright 1998-2001 Francois Desarmenien.
  6 # (c) Parse::Eyapp Copyright 2006 Casiano Rodriguez-Leon. Universidad de La Laguna.
  7 #   Don't edit this file, use source file "CalcSyntax.eypp" instead.
  8 #
  9 #           ANY CHANGE MADE HERE WILL BE LOST !
 10 #
 11 #####
 12 package Calc::Syntax;
```

La opción que recomiendo para documentar el módulo es escribir la documentación en un fichero aparte `Calc/Syntax.pod`.

El Programa Cliente

A continuación escribimos el programa cliente:

```
$ cd ../..
$ mkdir scripts
$ cd scripts/
$ vi usecalcsyntax.pl
$ cat -n usecalcsyntax.pl
  1 #!/usr/bin/perl -w -I../lib
  2 use strict;
  3 use Calc::Syntax;
  4 use Carp;
  5
  6 sub slurp_file {
  7     my $fn = shift;
  8     my $f;
  9
 10     local $/ = undef;
```

```

11  if (defined($fn)) {
12      open $f, $fn
13  }
14  else {
15      $f = \*STDIN;
16  }
17  my $input = <$f>;
18  return $input;
19 }
20
21 my $parser = Calc::Syntax->new();
22
23 my $input = slurp_file( shift() );
24 $parser->Run($input);

```

Ejecución

La ejecución muestra la *antiderivación a derechas construida por eyapp*:

```

$ cat prueba.exp
a=2*3

```

```

$ usecalcsyntax.pl prueba.exp
exp -> NUM (2)
exp -> NUM (3)
exp -> exp '*' exp
exp -> VAR '=' exp
line -> exp \n
input -> line *

```

Orden de Ejecución de las Acciones Semánticas

¿En que orden ejecuta YYParse las acciones? La respuesta es que el analizador generado por **eyapp** construye una derivación a derechas inversa y ejecuta las acciones asociadas a las reglas de producción que se han aplicado. Así, para la frase **2*3** la antiderivación es:

$$NUM + NUM \xleftarrow{NUM \leftarrow E} E + NUM \xleftarrow{NUM \leftarrow E} E + E \xleftarrow{E + E \leftarrow E} E$$

por tanto las acciones ejecutadas son las asociadas con las correspondientes reglas de producción:

1. La acción asociada con $E \rightarrow NUM$:

```
NUM          { print "exp -> NUM ($_[1])\n"; }
```

Esta instancia de **exp** tiene ahora como atributo 2 (pasado por el analizador léxico).

2. De nuevo la acción asociada con $E \rightarrow NUM$:

```
NUM          { print "exp -> NUM ($_[1])\n"; }
```

Esta nueva instancia de **exp** tiene como atributo 3.

3. La acción asociada con $E \rightarrow E * E$:

```
| exp '*' exp      { print "exp -> exp '*' exp\n"; }
```

Obsérvese que la antiderivación a derechas da lugar a un recorrido ascendente y a izquierdas del árbol:

$$E_3(E_1(NUM[2]), *, E_2(NUM[2]))$$

Los subíndices indican el orden de visita de los nodos/producciones.

8.3. Depuración de Errores

Las fuentes de error cuando se programa con una herramienta como **eyapp** son diversas. En esta sección trataremos tres tipos de error:

- Conflictos en la gramática: la gramática es ambigua o quizá - posiblemente debido a que en **eyapp** se usa un sólo símbolo de predicción - no queda claro para el generador que analizador sintáctico producir
- Ya no tenemos conflictos pero el analizador sintáctico generado no se comporta como esperábamos: no acepta frases correctas o construye un árbol erróneo para las mismas
- Ya hemos resuelto los conflictos y además el análisis sintáctico es correcto pero tenemos errores semánticos. Los errores se producen durante la ejecución de las acciones semánticas

Resolución de Conflictos

El siguiente programa **eyapp** contiene algunos errores. El lenguaje generado esta constituido por listas de Ds (por declaraciones) seguidas de listas de Ss (por sentencias) separadas por puntos y coma:

```
pl@nereida:~/LEyapp/examples$ cat -n Debug.eypp
1  %token D S
2
3  %{
4  our $VERSION = '0.01';
5  %}
6
7  %%
8  p:
9          ds ';' ss
10         | ss
11 ;
12
13 ds:
14         D ';' ds
15         | D
16         {
17             print "Reducing by rule:\n";
18             print "\tds -> D\n";
19             $_[1];
20         }
21 ;
22
23 ss:
24         S ';' ss
25         | S
26 ;
27
28 %%
29
30 my $tokenline = 0;
31
32 sub _Error {
33     my $parser = shift;
34     my ($token) = $parser->YYCurval;
35     my ($what) = $token ? "input: '$token'" : "end of input";
```

```

36         die "Syntax error near $what line num $tokenline\n";
37     }
38
39     my $input;
40
41     sub _Lexer {
42
43         for ($input) {
44             s{^\s}{}} and $tokenline += $1 =~ tr{\n}{};
45             return ('',undef) unless $_;
46             return ($1,$1) if s/^(.)//;
47         }
48         return ('',undef);
49     }
50
51     sub Run {
52         my ($self) = shift;
53
54         $input = shift;
55
56         return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error,
57                               yydebug => 0xF
58         );
59     }

```

Al compilar este programa con `eyapp` produce un mensaje de advertencia anunciándonos la existencia de un conflicto.

```

pl@nereida:~/LEyapp/examples$ eyapp Debug.eyp
1 shift/reduce conflict (see .output file)
State 4: shifts:
  to state    8 with ','

```

La existencia de advertencias da lugar a la creación de un fichero `Debug.output` conteniendo información sobre la gramática y el analizador construido.

Veamos los contenidos del fichero:

```

pl@nereida:~/LEyapp/examples$ cat -n Debug.output
 1 Warnings:
 2 -----
 3 1 shift/reduce conflict (see .output file)
 4 State 4: shifts:
 5   to state    8 with ','
 6
 7 Conflicts:
 8 -----
 9 State 4 contains 1 shift/reduce conflict
10
11 Rules:
12 -----
13 0:      $start -> p $end
14 1:      p -> ds ',' ss
15 2:      p -> ss
16 3:      ds -> D ',' ds
17 4:      ds -> D

```

```

18 5:      ss -> S ';' ss
19 6:      ss -> S
20
21 States:
22 -----
23 State 0:
24
25      $start -> . p $end      (Rule 0)
26
27      D      shift, and go to state 4
28      S      shift, and go to state 1
29
30      p      go to state 2
31      ss     go to state 3
32      ds     go to state 5
33
.. .....
55 State 4:
56
57      ds -> D . ';' ds      (Rule 3)
58      ds -> D .      (Rule 4)
59
60      ';'    shift, and go to state 8
61
62      ';'    [reduce using rule 4 (ds)]
63
.. .....
84 State 8:
85
86      ds -> D ';' . ds      (Rule 3)
87
88      D      shift, and go to state 4
89
90      ds     go to state 11
91
.. .....
112 State 12:
113
114      p -> ds ';' ss .      (Rule 1)
115
116      $default      reduce using rule 1 (p)
117
118
119 Summary:
120 -----
121 Number of rules      : 7
122 Number of terminals  : 4
123 Number of non-terminals : 4
124 Number of states    : 13

```

El problema según se nos anuncia ocurre en el estado 4. Como veremos mas adelante el analizador sintáctico generado por `Parse::Eyapp` es un autómata finito. Cada estado guarda información sobre las reglas que podrían aplicarse durante el análisis de la entrada. Veamos en detalle la información asociada con el estado 4:

```

55 State 4:
56
57     ds -> D . ';' ds      (Rule 3)
58     ds -> D .            (Rule 4)
59
60     ';'      shift, and go to state 8
61
62     ';'      [reduce using rule 4 (ds)]
63

```

Un estado es un conjunto de reglas de producción con un marcador en su parte derecha. La idea es que - si estamos en un estado dado - esas reglas son potenciales candidatos para la construcción del árbol sintáctico. Que lo sean o no dependerá de los terminales que se vean a continuación.

El punto que aparece en la parte derecha de una regla indica "posición" de lectura. Así el hecho de que en el estado cuatro aparezcan los items:

```

57     ds -> D . ';' ds      (Rule 3)
58     ds -> D .            (Rule 4)

```

Significa que si estamos en este estado es por que se leyó una D y se espera ver un punto y coma.

El comentario de la línea 60 indica que si el siguiente terminal es ; podríamos ir al estado 8. Obsérvese que el estado 8 contiene un item de la forma `ds -> D ';' . ds`. La marca recuerda ahora que se vió una D y un punto y coma. El comentario de la línea 62:

```

62     ';'      [reduce using rule 4 (ds)]

```

indica que `Parse::Eyapp` considera factible otro árbol cuando el token actual es punto y coma: Reducir por la regla `ds -> D`.

Para ilustrar el problema consideremos las frases `D;S` y `D;D;S`.

Para ambas frases, después de consumir la D el analizador irá al estado 4 y el terminal actual será punto y coma. Para la primera frase `D;S` la decisión correcta es utilizar ("reducir" en la jerga) la regla 4 `ds -> D`.

Para la segunda frase `D;D;S` la decisión correcta es conjeturar la regla 3 `ds -> D . ';' ds`. El analizador podría saber que regla es la correcta si consumiera el terminal que sigue al punto y coma: si es S se trata de la regla 4 y si es D de la regla 3. Los analizadores generados por Eyapp no "miran" mas allá del siguiente terminal y por tanto no están en condiciones de decidir.

Cualquier solución a este tipo de conflictos implica una reformulación de la gramática modificando prioridades o reorganizando las reglas. Reescribiendo la regla para `ds` recursiva por la derecha desaparece el conflicto:

```

l@nereida:~/LEyapp/examples$ sed -ne '/^ds:/,/~/p' Debug1.eyp | cat -n
1 ds:
2     ds ';' D
3     | D
4     {
5         print "Reducing by rule:\n";
6         print "\tds -> D\n";
7         $_[1];
8     }
9 ;

```

Ahora ante una frase de la forma `D ; ...` siempre hay que reducir por `ds -> D`. La antiderivación a derechas para `D;D;S` es:

Derivación	Árbol
D;D;S <= ds;D;S <= ds;S <= ds;ss <= p	p(ds(ds(D),',',D),',',ss(S))
o	

Mientras que la an-

tiderivación a derechas para D;S es:

Derivación	Árbol
D;S <= ds;S <= ds;ss <= p	p(ds(D),',',ss(S))

Recompilamos la nueva versión de la gramática. Las advertencias han desaparecido:

```
pl@nereida:~/LEyapp/examples$ eyapp Debug1.eyp
pl@nereida:~/LEyapp/examples$
```

Errores en la Construcción del Arbol

Escribimos el típico programa cliente:

```
pl@nereida:~/LEyapp/examples$ cat -n usedebug1.pl
 1  #!/usr/bin/perl -w
 2  # usetreebypass.pl prueba2.exp
 3  use strict;
 4  use Debug1;
 5
 6  sub slurp_file {
 7      my $fn = shift;
 8      my $f;
 9
10      local $/ = undef;
11      if (defined($fn)) {
12          open $f, $fn or die "Can't find file $fn!\n";
13      }
14      else {
15          $f = \*STDIN;
16      }
17      my $input = <$f>;
18      return $input;
19  }
20
21  my $input = slurp_file( shift() );
22
23  my $parser = Debug1->new();
24
25  $parser->Run($input);
```

y ejecutamos. Cuando damos la entrada D;S introduciendo algunos blancos y retornos de carro entre los terminales ocurre un mensaje de error:

```
casiano@cc111:~/LPLsrc/Eyapp/examples/debuggingtut$ usedebug1.pl
D
;
```



```

S
Reducing by rule:
    ds -> D
Syntax error near end of input line num 1

```

Como esta conducta es anómala activamos la opción `yydebug => 0xF` en la llamada a `YYParser`.

Es posible añadir un parámetro en la llamada a `YYParse` con nombre `yydebug` y valor el nivel de depuración requerido. Ello nos permite observar la conducta del analizador. Los posibles valores de depuración son:

Bit	Información de Depuración
0x01	Lectura de los terminales
0x02	Información sobre los estados
0x04	Acciones (shifts, reduces, accept ...)
0x08	Volcado de la pila
0x10	Recuperación de errores

Veamos que ocurre cuando damos la entrada `D;S` introduciendo algunos blancos y retornos de carro entre los terminales:

```

pl@nereida:~/LEyapp/examples$ usedebug1.pl
D

;

S
-----
In state 0:
Stack: [0]
Need token. Got >D<
Shift and go to state 4.
-----
In state 4:
Stack: [0,4]
Don't need token.
Reduce using rule 4 (ds --> D): Reducing by rule:
    ds -> D
Back to state 0, then go to state 5.
-----
In state 5:
Stack: [0,5]
Need token. Got ><
Syntax error near end of input line num 1

```

¿Que está pasando? Vemos que después de leer `D` el analizador sintáctico recibe un `end of file`. Algo va mal en las comunicaciones entre el analizador léxico y el sintáctico. Repasemos el analizador léxico:

```

pl@nereida:~/LEyapp/examples$ sed -ne '/sub.*_Lexer/,/^}/p' Debug1.eyy | cat -n
1  sub _Lexer {
2
3      for ($input) {

```

```

4         s{^\(s+)\}\} and $tokenline += $1 =~ tr{\n}{ };
5         return ('',undef) unless $_;
6         return ($1,$1) if s/^(.)//;
7     }
8     return ('',undef);
9 }

```

El error está en la línea 4. ¡Sólo se consume un blanco!. Escribimos una nueva versión Debug2.eyyp corrigiendo el problema:

```

pl@nereida:~/LEyapp/examples$ sed -ne '/sub.*_Lexer/,/^}/p' Debug2.eyyp | cat -n
1  sub _Lexer {
2
3      for ($input) {
4          s{^\(s+)\}\} and $tokenline += $1 =~ tr{\n}{ };
5          return ('',undef) unless $_;
6          return ($1,$1) if s/^(.)//;
7      }
8      return ('',undef);
9  }

```

Ahora el análisis parece funcionar correctamente:

```

pl@nereida:~/LEyapp/examples$ usedebug2.pl
D

```

;

S

```

-----
In state 0:
Stack:[0]
Need token. Got >D<
Shift and go to state 4.
-----
In state 4:
Stack:[0,4]
Don't need token.
Reduce using rule 4 (ds --> D): Reducing by rule:
    ds -> D
Back to state 0, then go to state 5.
-----
In state 5:
Stack:[0,5]
Need token. Got >;<
Shift and go to state 8.
-----
In state 8:
Stack:[0,5,8]
Need token. Got >S<
Shift and go to state 1.
-----
In state 1:
Stack:[0,5,8,1]

```

```

Need token. Got ><
Reduce using rule 6 (ss --> S): Back to state 8, then go to state 10.
-----
In state 10:
Stack:[0,5,8,10]
Don't need token.
Reduce using rule 1 (p --> ds ; ss): Back to state 0, then go to state 2.
-----
In state 2:
Stack:[0,2]
Shift and go to state 7.
-----
In state 7:
Stack:[0,2,7]
Don't need token.
Accept.

```

Errores en las Acciones Semánticas

Un tercer tipo de error ocurre cuando el código de una acción semántica no se conduce como esperamos.

Las acciones semánticas son convertidas en métodos anónimos del objeto analizador sintáctico generado por `Parse::Eyapp`. Puesto que son subrutinas anónimas no podemos utilizar comandos de ruptura del depurador como

```
b nombre # parar cuando se entre en la subrutina ''nombre''
```

o bien

```
c nombredesub # continuar hasta alcanzar la subrutina ''nombre''
```

Además el fichero cargado durante la ejecución es el `.pm` generado. El código en `Debug.pm` nos es ajeno - fue generado automáticamente por `Parse::Eyapp` - y puede resultar difícil encontrar en él las acciones semánticas que insertamos en nuestro esquema de traducción.

La estrategia a utilizar para poner un punto de parada en una acción semántica se ilustra en la siguiente sesión con el depurador. Primero arrancamos el depurador con el programa y usamos la opción `f file` para indicar el nuevo "fichero de vistas" por defecto:

```
pl@nereida:~/LEyapp/examples$ perl -wd usedebug2.pl
```

```
Loading DB routines from perl5db.pl version 1.28
Editor support available.
```

```
Enter h or 'h h' for help, or 'man perldebug' for more help.
```

```

main:.(usedebug2.pl:21):          my $input = slurp_file( shift() );
  DB<1> f Debug2.eyp
1          2          #line 3 "Debug2.eyp"
3
4:          our $VERSION = '0.01';
5
6          7          8          9          10

```

Ahora usamos la orden `b 18` para poner un punto de ruptura en la línea 18. El comando `l` muestra las correspondientes líneas del fichero `.eyp`:

```

DB<2> b 18
DB<3> l
11      12      13      14      15      16      #line 17 "Debug2.eyp"
17
18:b      print "Reducing by rule:\n";
19:      print "\tds -> D\n";
20:      $_[1];

```

La orden c (continuar) hace que el programa se ejecute hasta alcanzar el único punto de ruptura en la línea 18 de Debug2.eyp:

```

DB<3> c
D

```

```
;
```

```
S
```

```

Debug2::CODE(0x85129d8)(Debug2.eyp:18):
18:      print "Reducing by rule:\n";
DB<3> n
Reducing by rule:

```

En este momento podemos usar cualesquiera comandos del depurador para visualizar el estado interno de nuestro programa y determinar la causa de una potencial conducta anómala de la acción semántica:

```

Debug2::CODE(0x85129d8)(Debug2.eyp:19):
19:      print "\tds -> D\n";
DB<3> x $_[0]{GRAMMAR}
0  ARRAY(0x8538360)
  0  ARRAY(0x855aa88)
    0  '_SUPERSTART'
    1  '$start'
    2  ARRAY(0x855ab60)
      0  'p'
      1  '$end'
    3  0
  1  ARRAY(0x855a890)
    0  'p_1'
    1  'p'
    2  ARRAY(0x855a8fc)
      0  'ds'
      1  ';'
      2  'ss'
    3  0
  2  ARRAY(0x855a800)
    0  'p_2'
    1  'p'
    2  ARRAY(0x855a830)
      0  'ss'
    3  0
  3  ARRAY(0x855a764)
    0  'ds_3'
    1  'ds'
    2  ARRAY(0x855a7a0)

```

```

        0  'ds'
        1  ';'
        2  'D'
    3  0
4  ARRAY(0x85421d4)
    0  'ds_4'
    1  'ds'
    2  ARRAY(0x855a6e0)
        0  'D'
    3  0
5  ARRAY(0x8538474)
    0  'ss_5'
    1  'ss'
    2  ARRAY(0x854f9c8)
        0  'S'
        1  ';'
        2  'ss'
    3  0
6  ARRAY(0x85383b4)
    0  'ss_6'
    1  'ss'
    2  ARRAY(0x85383f0)
        0  'S'
    3  0
DB<4>

```

Con el comando `c` podemos hacer que la ejecución continúe, esta vez hasta a finalización del programa

```

DB<3> c
Debugged program terminated. Use q to quit or R to restart,
use o inhibit_exit to avoid stopping after program termination,
h q, h R or h o to get additional info.
DB<3>

```

8.4. Acciones y Acciones por Defecto

En el ejemplo anterior la acción semántica asociada con cada una de las reglas de producción es básicamente la misma: escribir la regla. Parece lógico intentar factorizar todo ese código común.

Las Acciones en eyapp

En `Parse::Eyapp` las acciones semánticas son convertidas en subrutinas anónimas. Mas bien en métodos anónimos dado que el primer argumento (`$_[0]`) de toda acción semántica es una referencia al objeto analizador sintáctico. Los restantes parámetros se corresponden con los *atributos de los símbolos* en la parte derecha de la regla de producción (`$_[1] ...`). Esto es, si la regla es:

$$A \rightarrow X_1 X_2 \dots X_n \quad \{ \text{action}(@_) \}$$

cuando se "reduzca" por la regla $A \rightarrow X_1 X_2 \dots X_n$ la acción será llamada con argumentos

```
action(parserref, $X_1, $X_2, ..., $X_n)
```

donde `$X_1`, `$X_2`, etc. denotan los atributos de X_1 , X_2 , etc. y `parserref` es una referencia al objeto analizador sintáctico.

El valor retornado por la acción semántica es asignado como valor del atributo del lado izquierdo de la regla A .

Cuando una regla $A \rightarrow X_1X_2 \dots X_n$ no tiene acción semántica asociada se ejecuta la *acción por defecto*. La acción por defecto en `eyapp` es `{ $_[1] }`:

$$A \rightarrow X_1X_2 \dots X_n \quad \{ \$_[1] \}$$

El efecto de tal acción es asignar al atributo de la variable sintáctica en la parte izquierda el atributo del primer símbolo en la parte derecha. Si la parte derecha fuera vacía se asigna `undef`.

La Directiva `%defaultaction`

La directiva `%defaultaction` permite especificar una acción por defecto alternativa. Observe esta nueva versión del programa anterior:

```
pl@nereida:~/LEyapp/examples$ cat -n CalcSyntaxDefaultWithHOPLexer.eypp
1  # CalcSyntaxDefaultWithHOPLexer.eypp
2  %right  '='
3  %left   '- ' '+'
4  %left   '* ' '/'
5  %left   NEG
6  %right  '^'
7
8  %defaultaction {
9      my $parser = shift;
10
11      print $parser->YYLhs." --> ";
12      my $i = 0;
13      for ($parser->YYRightside) {
14          print "$_";
15          print "[$_[ $i]]" if /NUM|VAR/;
16          print " ";
17          $i++;
18      }
19      print "\n";
20 }
21
22 %{
23 use HOP::Lexer qw(string_lexer);
24 %}
25 %%
26 line: (exp '\n')*
27 ;
28
29 exp:
30     NUM
31     | VAR
32     | VAR '=' exp
33     | exp '+' exp
34     | exp '- ' exp
35     | exp '* ' exp
36     | exp '/' exp
37     | '- ' exp %prec NEG
38     | exp '^' exp
39     | '(' exp ')'
40 ;
```

El método `YYLhs` (línea 11) retorna el nombre de la variable sintáctica asociada con la regla actual. El método `YYRightside` (línea 13) devuelve la lista de nombres de los símbolos de la parte derecha de la regla de producción actual.

Observe como la modificación de la acción con defecto nos permite eliminar todas las acciones del cuerpo de la gramática.

Al ejecutar el cliente obtenemos una salida similar a la del programa anterior:

```
pl@nereida:~/LEyapp/examples$ eyapp CalcSyntaxDefaultWithHOPlexer.eyp
pl@nereida:~/LEyapp/examples$ cat -n prueba4.exp
 1  4*3
pl@nereida:~/LEyapp/examples$ usecalcsyntaxdefaultwithhoplexer.pl prueba4.exp | cat -n
 1  exp --> NUM[4]
 2  exp --> NUM[3]
 3  exp --> exp * exp
 4  line --> STAR-2
```

La Opción `-v` de `eyapp`

Posiblemente llame su atención la línea `line --> STAR-2` en la salida. Es consecuencia de la existencia de la línea 26 que contiene una aplicación del operador de repetición `*` a una expresión entre paréntesis:

```
26 line: (exp '\n')*
```

De hecho la gramática anterior es desplegada como se indica en el fichero `CalcSyntaxDefaultWithHOPlexer.output` obtenido al compilar usando la opción `-v` de `eyapp`:

```
pl@nereida:~/LEyapp/examples$ eyapp -v CalcSyntaxDefaultWithHOPlexer.eyp
pl@nereida:~/LEyapp/examples$ ls -ltr | tail -1
-rw-r--r-- 1 pl users 8363 2007-10-30 11:45 CalcSyntaxDefaultWithHOPlexer.output
pl@nereida:~/LEyapp/examples$ sed -ne '/Rules/,/^$/p' CalcSyntaxDefaultWithHOPlexer.output
Rules:
-----
0:      $start -> line $end
1:      PAREN-1 -> exp '\n'
2:      STAR-2 -> STAR-2 PAREN-1
3:      STAR-2 -> /* empty */
4:      line -> STAR-2
5:      exp -> NUM
6:      exp -> VAR
7:      exp -> VAR '=' exp
8:      exp -> exp '+' exp
9:      exp -> exp '-' exp
10:     exp -> exp '*' exp
11:     exp -> exp '/' exp
12:     exp -> '-' exp
13:     exp -> exp '^' exp
14:     exp -> '(' exp ')'
```

Análisis Léxico con `HOP::Lexer`

En esta ocasión hemos elegido un método alternativo para construir el analizador léxico: Usar el módulo `HOP::Lexer`.

El módulo `HOP::Lexer` es cargado en las líneas 22-24. Los delimitadores `%{` y `%}` permiten introducir código Perl en cualquier lugar de la sección de cabecera.

Estudie la documentación del módulo `HOP::Lexer` y compare la solución que sigue con la dada en la sección anterior (véase 8.2). Este es el código en la sección de cola:

```

42 %%
43
44 sub _Error {
45     .....
46 }
47
48 sub Run {
49     my($self)=shift;
50     my $input = shift;
51
52     my @lexemes = (
53         ['SPACE', qr{[ \t]+}, sub { () } ],
54         ['NUM',   qr{[0-9]+(?:\.[0-9]+)?}],
55         ['VAR',   qr{[A-Za-z][A-Za-z0-9_]*}],
56         ["\n",    qr{\n}],
57         ['ANY',   qr{.}],
58         sub {
59             my ($label, $val) = @_;
60             return [ $val, $val ];
61         }
62     ]
63 );
64
65 my $lexer = string_lexer($input, @lexemes);
66
67 return $self->YYParse(
68     yylex => sub {
69         my $parser = shift;
70         my $token = $lexer->();
71         return @$token if defined($token);
72         return ('', undef);
73     },
74     yyerror => \&_Error,
75     #yydebug => 0xF
76 );
77 }

```

La Opción yydebug

Cuando YYParse se llama con la opción yydebug activada (descomente la línea 81 si quiere ver el efecto) la ejecución del analizador produce un informe detallado (version de eyapp 0.83 o superior) de su avance durante el análisis:

```

pl@nereida:~/LEyapp/examples$ usecalcsyntaxdefaultwithhoplexer.pl prueba4.exp >& warn
pl@nereida:~/LEyapp/examples$ head -30 warn

```

```

-----
In state 0:
Stack:[0]
Don't need token.
Reduce using rule 3 (STAR-2 --> /* empty */): Back to state 0, then go to state 1.
-----

```

```

In state 1:
Stack:[0,1]
Need token. Got >NUM<
Shift and go to state 4.

```



```

-----
In state 4:
Stack: [0,1,4]
Don't need token.
Reduce using rule 5 (exp --> NUM): Back to state 1, then go to state 5.
-----

In state 5:
Stack: [0,1,5]
Need token. Got >*<
Shift and go to state 13.
-----

In state 13:
Stack: [0,1,5,13]
Need token. Got >NUM<
Shift and go to state 4.
-----

In state 4:
Stack: [0,1,5,13,4]
Don't need token.
Reduce using rule 5 (exp --> NUM): Back to state 13, then go to state 21.

```

8.5. Traducción de Infijo a Postfijo

Las acción por defecto (en yacc `$$ = $1`) puede ser modificada de manera parecida a como lo hace la variable `RD_AUTOACTION` en `Parse::RecDescent` usando la directiva `%defaultaction`:

```
%defaultaction { perl code }
```

El código especificado en la acción por defecto se ejecutará siempre que se produzca una reducción por una regla en la que no se haya explicitado código.

El siguiente ejemplo traduce una expresión en infijo como `a=b*-3` a una en postfijo como `a b 3 NEG * = :`

```

# File Postfix.y
%right  '='
%left   '-' '+'
%left   '*' '/'
%left   NEG
%{
use IO::Interactive qw(interactive);
%}

%defaultaction {
    my $self = shift;

    return "$_[0]" if (@_==1); # NUM and VAR
    return "$_[0] $_[2] $_[1]" if (@_==3); # other operations
    die "Fatal Error\n";
}

%%
line: (exp { print "$_[1]\n" })?
;

exp:      NUM

```

```

|   VAR
|   VAR '=' exp
|   exp '+' exp
|   exp '-' exp
|   exp '*' exp
|   exp '/' exp
|   '-' exp %prec NEG { "$_[2] NEG" }
|   '(' exp ')' { $_[2] }
;

%%

sub _Error {
    my($token)=$_[0]->YYCurval;
    my($what)= $token ? "input: '$token'" : "end of input";
    my @expected = $_[0]->YYExpect();

    local $" = ', ';
    die "Syntax error near $what. Expected one of these tokens: @expected\n";
}

my $x;

sub _Lexer {
    my($parser)=shift;

    for ($x) {
        s/^\s+//;
        $_ eq '' and return('','undef');

        s/^\([0-9]+(?:\.[0-9]+)?\)// and return('NUM',$1);
        s/^\([A-Za-z][A-Za-z0-9_]*\)// and return('VAR',$1);
        s/^(.)//s and return($1,$1);
    }
}

sub Run {
    my($self)=shift;
    print {interactive} "Write an expression: ";
    $x = <>;
    $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error,
        #yydebug => 0xFF
    );
}

```

8.6. Práctica: Traductor de Términos a GraphViz

El objetivo de esta práctica es el diseño e implementación de un lenguaje de domino específico (DSL) usando Parse::Eyapp.

- Escriba una gramática que reconozca el lenguaje de los términos que describen árboles. Frases como

```
MULT(NUM,ADD(NUM,NUM))
```

o

```
MULT\(*\) (NUM\ (2\) , ADD\ (+\) (NUM\ (3\) , NUM\ (4\)))
```

deberían ser aceptadas.

Un término no es mas que una secuencia de nombres de nodo seguida de un paréntesis abrir y una lista de términos separados por comas, terminada por el correspondiente paréntesis cerrar. Los nombres de nodo admiten cualesquiera caracteres que no sean los paréntesis. Sin embargo, los nombres de nodo pueden contener paréntesis escapados.

- Instale GraphViz y el módulo GraphViz de Leon Brocard en su ordenador. Con este módulo es posible generar imágenes a partir de una descripción de alto nivel de un grafo. Por ejemplo, el programa Perl:

```
pl@nereida:~/Lbook$ cat -n ast.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3  use GraphViz;
 4
 5  my $g = GraphViz->new();
 6
 7  $g->add_node('MULT(*)', shape => 'box', color => 'blue');
 8  $g->add_node('NUM(2)', shape => 'box', color => 'blue');
 9  $g->add_node("ADD(+)", shape => 'box', color => 'blue');
10  $g->add_node("NUM(3)", shape => 'box', color => 'blue');
11  $g->add_node("NUM(4)", shape => 'box', color => 'blue');
12
13  $g->add_edge('MULT(*)' => "NUM(2)", label => 'left', color => 'blue');
14  $g->add_edge('MULT(*)' => "ADD(+)", label => 'right', color => 'blue');
15  $g->add_edge("ADD(+)" => "NUM(3)", label => 'left', color => 'blue');
16  $g->add_edge("ADD(+)" => "NUM(4)", label => 'right', color => 'blue');
17
18  open my $f, ">", "perlexamples/ast234.png";
19  print $f $g->as_png;
20  close($f);
```

Genera un fichero **png** (es posible generar otros formatos, si se prefiere) conteniendo el gráfico del árbol:

- Extienda el reconocedor de términos árbol para que genere como traducción un programa Perl que genera la correspondiente imagen **.png**. Sigue un ejemplo de llamada:

```
$ tree2png -t 'MULT(NUM,ADD(NUM,NUM))' -o ast.png
```

- Si lo prefiere, puede usar el lenguaje **dot** como lenguaje objeto. El lenguaje **dot** permite generar gráficos a partir de una descripción textual de un grafo.

Por ejemplo, a partir del programa **dot**:

```
~/src/perl/graphviz$ cat -n hierarchy.dot
 1  digraph hierarchy {
 2
 3  nodesep=1.0 // increases the separation between nodes
 4
 5  node [color=Red,fontname=Courier]
```

```

6  edge [color=Blue, style=dashed] //setup options
7
8  Boss->{ John Jack } // the boss has two employees
9
10 {rank=same; John Jack} //they have the same rank
11
12 John -> Al // John has a subordinate
13
14 John->Jack [dir=both] // but is still on the same level as Jack
15 }

```

mediante el comando:

```
~/src/perl/graphviz$ dot hierarchy.dot -Tpng > h.png
```

producimos la siguiente imagen:

8.7. Práctica: Gramática Simple en Parse::Eyapp

Escriba usando Parse::Eyapp un analizador sintáctico PL::Simple para la siguiente gramática:

0	$p \rightarrow ds\ ss$
1	$p \rightarrow ss$
2	$ds \rightarrow d\ ';\ ' ds$
3	$ds \rightarrow d\ ';\ '$
4	$d \rightarrow INT\ il$
5	$d \rightarrow STRING\ il$
6	$ss \rightarrow s\ ';\ ' ss$
7	$ss \rightarrow s$
8	$s \rightarrow ID\ '='\ e$
9	$s \rightarrow PRINT\ nel$
10	$s \rightarrow \epsilon$
11	$e \rightarrow e\ '+'\ t$
12	$e \rightarrow e\ '-'\ t$
13	$e \rightarrow t$
14	$t \rightarrow t\ '*' \ f$
15	$t \rightarrow t\ '/' \ f$
16	$t \rightarrow f$
17	$f \rightarrow '('\ e\ ')'$
18	$f \rightarrow ID$
19	$f \rightarrow NUM$
20	$f \rightarrow STR$
21	$f \rightarrow FUNCTIONID\ '('\ el\ ')'$
22	$il \rightarrow ID\ ','\ il$
23	$il \rightarrow ID$
24	$el \rightarrow nel$
25	$el \rightarrow \epsilon$
26	$nel \rightarrow e\ ','\ nel$
27	$nel \rightarrow e$

El terminal ID corresponde a identificador y NUM a un número entero sin signo. El terminal STR a cadenas de dobles comillas. El terminal FUNCTIONID viene dado por la expresión regular `min|max|concat`. Declare dichos terminales en la cabecera usando la directiva `%token` (lea la documentación de Parse::Eyapp). Una directiva como:

%token ID FUNCTIONID NUM

declara los terminales pero no les da prioridad.

- Describa mediante ejemplos el lenguaje generado por esta gramática. ¿Es ambigua la gramática?
- Admita comentarios como en C.
- Almacene los números de línea asociados con los terminales y utilice esa información en los diagnósticos de error.
- Para la escritura de las reglas del forma $s \rightarrow \epsilon$ repase los consejos dados en el párrafo 8.2
- Expanda la gramática para que incluya operaciones lógicas ($a < b$, $a == b$, etc.). Observe que tales expresiones tiene menor prioridad que los operadores aditivos: $a = b < c + 2$ es interpretada como $a = b < (c + 2)$
- Escriba la documentación en un fichero de documentación aparte `PL/Simple.pod`.
- Escriba las pruebas. En las pruebas asegúrese de poner tanto entradas correctas como erróneas.
- Adjunte un estudio de cubrimiento con su práctica. Añada los ficheros con el informe al **MANIFEST**. Ubíquelos en un directorio de informes.

8.8. El Método YYName y la Directiva %name

Construyendo una Representación del Árbol

En `eyapp` toda regla de producción tiene un nombre. El nombre de una regla puede ser dado explícitamente por el programador mediante la directiva `%name`. Por ejemplo, en el fragmento de código que sigue se da el nombre `ASSIGN` a la regla `exp: VAR '=' exp`. Si no se da un nombre explícito, la regla tendrá un nombre implícito. El nombre implícito de una regla se forma concatenando el nombre de la variable sintáctica en la parte izquierda con un subguión y el número de orden de la regla de producción: `Lhs_#`. Evite dar nombres con ese patrón a sus reglas de producción. Los patrones de la forma `/${lhs}_\d+$/` donde `/${lhs}` es el nombre de la variable sintáctica están reservados por `eyapp`.

```
pl@nereida:~/LEyapp/examples$ cat -n Lhs.eyp
1  # Lhs.eyp
2
3  %right  '='
4  %left   '- ' '+'
5  %left   '* ' '/'
6  %left   NEG
7
8  %defaultaction {
9      my $self = shift;
10     my $name = $self->YYName();
11     bless { children => [ grep {ref($_)} @_ ] }, $name;
12 }
13
14 %%
15 input:
16         /* empty */
17         { [] }
18     |   input line
19         {
20         push @{$_[1]}, $_[2] if defined($_[2]);
```

```

21             $_[1]
22         }
23     ;
24
25 line:      '\n'      { }
26         | exp '\n'    { $_[1] }
27     ;
28
29 exp:
30     NUM    { $_[1] }
31     | VAR  { $_[1] }
32     | %name ASSIGN
33     VAR '=' exp
34     | %name PLUS
35     exp '+' exp
36     | %name MINUS
37     exp '-' exp
38     | %name TIMES
39     exp '*' exp
40     | %name DIV
41     exp '/' exp
42     | %name UMINUS
43     '-' exp %prec NEG
44     | '(' exp ')' { $_[2] }
45 ;

```

Dentro de una acción semántica el nombre de la regla actual puede ser recuperado mediante el método `YYName` del analizador sintáctico.

La acción por defecto (líneas 8-12) computa como atributo de la parte izquierda una referencia a un objeto que tiene como clase el nombre de la regla. Ese objeto tiene un atributo `children` que es una referencia a la lista de hijos del nodo. La llamada a `grep`

```

11     bless { children => [ grep {ref($_)} @_ ] }, $name;

```

tiene como efecto obviar (no incluir) aquellos hijos que no sean referencias. Observe que el analizador léxico sólo retorna referencias a objetos para los terminales `NUM` y `VAR`.

La Cola

```

47 %%
48
49 sub _Error {
50     exists $_[0]->YYData->{ERRMSG}
51     and do {
52         print $_[0]->YYData->{ERRMSG};
53         delete $_[0]->YYData->{ERRMSG};
54         return;
55     };
56     print "Syntax error.\n";
57 }
58
59 sub _Lexer {
60     my($parser)=shift;
61
62     for ($parser->YYData->{INPUT}) {

```

```

63         s/^[ \t]+//;
64         return('',undef) unless $_;
65         s/^[0-9]+(?:\.[0-9]+)?//
66             and return('NUM', bless { attr => $1}, 'NUM');
67         s/^[A-Za-z][A-Za-z0-9_]*//
68             and return('VAR',bless {attr => $1}, 'VAR');
69         s/^(.)/s
70             and return($1, $1);
71     }
72     return('',undef);
73 }
74
75 sub Run {
76     my($self)=shift;
77
78     $self->YYData->{INPUT} = <>;
79     return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
80 }

```

Para los terminales NUM y VAR el analizador léxico devuelve una referencia a un objeto. No así para el resto de los terminales.

El Programa Cliente

```

pl@nereida:~/LEyapp/examples$ cat -n uselhs.pl
 1  #!/usr/bin/perl -w
 2  use Lhs;
 3  use Data::Dumper;
 4
 5  $parser = new Lhs();
 6  my $tree = $parser->Run;
 7  $Data::Dumper::Indent = 1;
 8  if (defined($tree)) { print Dumper($tree); }
 9  else { print "Cadena no válida\n"; }

```

Ejecución

Al darle la entrada $a=(2+3)*b$ el analizador produce el árbol

```
ASSIGN(TIMES(PLUS(NUM[2],NUM[3]), VAR[b]))
```

Veamos el resultado de una ejecución:

```

pl@nereida:~/LEyapp/examples$ uselhs.pl
a=(2+3)*b
$VAR1 = [
  bless( {
    'children' => [
      bless( { 'attr' => 'a' }, 'VAR' ),
      bless( {
        'children' => [
          bless( {
            'children' => [
              bless( { 'attr' => '2' }, 'NUM' ),
              bless( { 'attr' => '3' }, 'NUM' )

```

```

    ]
    }, 'PLUS' ),
    bless( { 'attr' => 'b' }, 'VAR' )
  ]
}, 'TIMES' )
]
}, 'ASSIGN' )
];

```

Cambiando el Nombre de la Regla Dinámicamente

Es posible cambiar en tiempo de ejecución el nombre de una regla. Observe la siguiente variante del ejemplo anterior en el que se cambia - en tiempo de ejecución - el nombre de la regla del menos binario.

```

29  exp:
30      NUM    { $_[1] }
31      |     VAR    { $_[1] }
32      |     %name ASSIGN
33      VAR '=' exp
34      |     %name PLUS
35      exp '+' exp
36      |     %name MINUS
37      exp '-' exp
38      {
39          my $self = shift;
40          $self->YYName('SUBTRACT'); # rename it
41          $self->YYBuildAST(@_); # build the node
42      }
43      |     %name TIMES
44      exp '*' exp
45      |     %name DIV
46      exp '/' exp
47      |     %name UMINUS
48      '-' exp %prec NEG
49      |     '(' exp ')' { $_[2] }
50  ;

```

Al ejecutar el cliente observamos como los nodos han sido bendecidos en la clase SUBTRACT:

```

pl@nereida:~/LEyapp/examples$ useyynamedynamic.pl
2-b
$VAR1 = [
  bless( {
    'children' => [
      bless( {
        'attr' => '2'
      }, 'NUM' ),
      bless( {
        'attr' => 'b'
      }, 'VAR' )
    ]
  }, 'SUBTRACT' )
];

```


8.9. Construyendo el Arbol de Análisis Sintactico Mediante Directivas

Directivas para la Construcción del Arbol Sintáctico

La directiva `%tree` hace que `Parse::Eyapp` genere una estructura de datos que representa al árbol sintáctico. Haciendo uso de las directivas adecuadas podemos controlar la forma del árbol.

Los nodos del árbol sintáctico generado por la producción $A \rightarrow X_1 \dots X_n$ son objetos (de tipo hash) bendecidos en una clase con nombre `A_#`, esto es, el de la variable sintáctica en el lado izquierdo seguida de un número de orden. Los nodos retornados por el analizador léxico son bendecidos en la clase especial `TERMINAL`. Los nodos tienen además un atributo `children` que referencia a la lista de nodos hijos $X_1 \dots X_n$ del nodo.

La directiva `%name CLASSNAME` permite modificar el nombre por defecto de la clase del nodo para que sea `CLASSNAME`. Por ejemplo, cuando es aplicada en el siguiente fragmento de código:

```
25  exp:
..      .....
32      | %name PLUS
33      exp '+' exp
```

Hace que el nodo asociado con la regla $exp \rightarrow exp + exp$ pertenezca a la clase `PLUS` en vez de a una clase con un nombre poco significativo como `exp_25`:

```
pl@nereida:~/LEyapp/examples$ cat -n CalcSyntaxTree.eyp
```

```
1  # CalcSyntaxTree.eyp
2  %right '='
3  %left  '-' '+'
4  %left  '*' '/'
5  %left  NEG
6  %right '^'
7
8  %tree
9
10 %{\
11 sub TERMINAL::info {
12   my $self = shift;
13
14   $self->attr;
15 }
16
17 %}\
18 %%
19
20 line:
21     %name EXP
22     exp '\n'
23 ;
24
25 exp:
26     %name NUM
27     NUM
28     | %name VAR
29     VAR
30     | %name ASSIGN
31     VAR '=' exp
```

```

32  | %name PLUS
33      exp '+' exp
34  | %name MINUS
35      exp '-' exp
36  | %name TIMES
37      exp '*' exp
38  | %name DIV
39      exp '/' exp
40  | %name UMINUS
41      '-' exp %prec NEG
42  | %name EXP
43      exp '^' exp
44  | %name PAREN
45      '(' exp ')'
46  ;
47
48  %%
.. # Exactamente igual que en el ejemplo anterior

```

La forma que tiene el árbol construido mediante la directiva `%tree` puede ser modificada. Por defecto, todos aquellos terminales que aparecen en definidos en el programa **eyapp** mediante el uso de apóstrofes son eliminados del árbol. En la jerga "Eyapp" un *token sintáctico* es uno que será eliminado (podado) del árbol sintáctico. Por contra, un *token semántico* es uno que aparecerá en el árbol. Por defecto los tokens definidos simbólicamente como `NUM` o `VAR` son "semánticos" y aparecerán como un nodo de tipo `TERMINAL` en el árbol de análisis sintáctico. Estos estatus por defecto pueden ser modificados mediante las directivas `%semantic token` y `%syntactic token`.

Los token sintácticos no forman parte del árbol construido. Así pues, en el ejemplo que nos ocupa, los terminales `'='`, `'-'`, `'+'`, `'*'` y `'/'` serán, por defecto, eliminados del árbol sintáctico.

El Cliente

```

pl@nereida:~/LEyapp/examples$ cat -n usecalcsyntaxtree.pl
1  #!/usr/bin/perl -w
2  # usecalcsyntaxtree.pl prueba2.exp
3  use strict;
4  use CalcSyntaxTree;
5
6  sub slurp_file {
7      my $fn = shift;
8      my $f;
9
10     local $/ = undef;
11     if (defined($fn)) {
12         open $f, $fn
13     }
14     else {
15         $f = \*STDIN;
16     }
17     my $input = <$f>;
18     return $input;
19 }
20
21 my $parser = CalcSyntaxTree->new();
22

```

```

23 my $input = slurp_file( shift() );
24 my $tree = $parser->Run($input);
25
26 $Parse::Eyapp::Node::INDENT = 2;
27 print $tree->str."\n";

```

El método `str` de `Parse::Eyapp::Node` devuelve una cadena describiendo el árbol de análisis sintáctico enraizado en el nodo que se ha pasado como argumento.

El método `str` cuando visita un nodo comprueba la existencia de un método `info` para la clase del nodo. Si es así el método será llamado. Obsérvese como en la cabecera del programa `Eyapp` proveemos un método `info` para los nodos `TERMINAL`.

```

8  %tree
9
10 %{
11 sub TERMINAL::info {
12     my $self = shift;
13
14     $self->attr;
15 }
16
17 %}
18 %%

```

Los nodos `TERMINAL` del árbol sintáctico tienen un atributo `attr` que guarda el valor pasado para ese terminal por el analizador léxico.

La variable de paquete `$Parse::Eyapp::Node::INDENT` controla el formato de presentación usado por `Parse::Eyapp::Node::str` : Si es 2 cada paréntesis cerrar que este a una distancia mayor de `$Parse::Eyapp::Node::LINESEP` líneas será comentado con el tipo del nodo.

Ejecución

```

pl@nereida:~/LEyapp/examples$ cat prueba2.exp
a=(2+b)*3
pl@nereida:~/LEyapp/examples$ usecalcsyntaxtree.pl prueba2.exp | cat -n
 1
 2 EXP(
 3     ASSIGN(
 4         TERMINAL[a] ,
 5         TIMES(
 6             PAREN(
 7                 PLUS(
 8                     NUM(
 9                         TERMINAL[2]
10                     ),
11                     VAR(
12                         TERMINAL[b]
13                     )
14                 ) # PLUS
15             ) # PAREN,
16             NUM(
17                 TERMINAL[3]
18             )
19         ) # TIMES
20     ) # ASSIGN
21 ) # EXP

```

Estructura Interna de los Árboles Sintácticos

Para entender mejor la representación interna que `Parse::Eyapp` hace del árbol ejecutemos de nuevo el programa con la ayuda del depurador:

```
pl@nereida:~/LEyapp/examples$ perl -wd usecalcsyntaxtree.pl prueba2.exp
main::(usecalcsyntaxtree.pl:21):          my $parser = CalcSyntaxTree->new();
  DB<1> l 21-27    # Listamos las líneas de la 1 a la 27
21==>  my $parser = CalcSyntaxTree->new();
22
23:    my $input = slurp_file( shift() );
24:    my $tree = $parser->Run($input);
25
26:    $Parse::Eyapp::Node::INDENT = 2;
27:    print $tree->str."\n";
  DB<2> c 27    # Continuamos la ejecución hasta alcanzar la línea 27
main::(usecalcsyntaxtree.pl:27):          print $tree->str."\n";
  DB<3> x $tree
0 EXP=HASH(0x83dec7c)    # El nodo raíz pertenece a la clase EXP
  'children' => ARRAY(0x83df12c) # El atributo 'children' es una referencia a un array
    0 ASSIGN=HASH(0x83decdc)
      'children' => ARRAY(0x83df114)
        0 TERMINAL=HASH(0x83dec40) # Nodo construido para un terminal
          'attr' => 'a' # Atributo del terminal
          'children' => ARRAY(0x83dee38) # Para un terminal debera ser una lista
            empty array                # vacía
          'token' => 'VAR'
        1 TIMES=HASH(0x83df084)
          'children' => ARRAY(0x83df078)
            0 PAREN=HASH(0x83ded9c)
              'children' => ARRAY(0x83defac)
                0 PLUS=HASH(0x83dee74)
                  'children' => ARRAY(0x83deef8)
                    0 NUM=HASH(0x83dedc0)
                      'children' => ARRAY(0x832df14)
                        0 TERMINAL=HASH(0x83dedfc)
                          'attr' => 2
                          'children' => ARRAY(0x83dedd8)
                            empty array
                          'token' => 'NUM'
                        1 VAR=HASH(0x83dec0)
                          'children' => ARRAY(0x83dee2c)
                            0 TERMINAL=HASH(0x83deec8)
                              'attr' => 'b'
                              'children' => ARRAY(0x83dee98)
                                empty array
                              'token' => 'VAR'
                            1 NUM=HASH(0x83dee44)
                              'children' => ARRAY(0x83df054)
                                0 TERMINAL=HASH(0x83df048)
                                  'attr' => 3
                                  'children' => ARRAY(0x83dece8)
                                    empty array
                                  'token' => 'NUM'
```

Observe que un nodo es un objeto implantado mediante un hash. El objeto es bendecido en la clase/paquete que se especifico mediante la directiva `%name`. Todas estas clases heredan de la clase `Parse::Eyapp::Node` y por tanto disponen de los métodos proveídos por esta clase (en particular el método `str` usado en el ejemplo). Los nodos disponen de un atributo `children` que es una referencia a la lista de nodos hijo. Los nodos de la clase `TERMINAL` son construidos a partir de la pareja (`token`, `atributo`) proveída por el analizador léxico. Estos nodos disponen además del atributo `attr` que encapsula el atributo del terminal.

Terminales Semánticos y Terminales Sintácticos

Observe como en el árbol anterior no existen nodos `TERMINAL[+]` ni `TERMINAL[*]` ya que estos terminales fueron definidos mediante apóstrofes y son - por defecto - terminales sintácticos.

Ejercicio 8.9.1. *Modifique el programa eyapp anterior para que contenga la nueva línea 2:*

```
pl@nereida:~/LEyapp/examples$ head -9 CalcSyntaxTree3.eyp | cat -n
 1 # CalcSyntaxTree.eyp
 2 %semantic token '=' '-' '+' '*' '/' '^'
 3 %right '='
 4 %left '-' '+'
 5 %left '*' '/'
 6 %left NEG
 7 %right '^'
 8
 9 %tree
```

*Escriba un programa cliente y explique el árbol que resulta para la entrada `a=(2+b)*3`:*

```
pl@nereida:~/LEyapp/examples$ cat prueba2.exp
a=(2+b)*3
pl@nereida:~/LEyapp/examples$ usecalcsyntaxtree3.pl prueba2.exp | cat -n
 1
 2 EXP(
 3   ASSIGN(
 4     TERMINAL[a],
 5     TERMINAL[=],
 6     TIMES(
 7       PAREN(
 8         PLUS(
 9         NUM(
10           TERMINAL[2]
11         ),
12         TERMINAL[+],
13         VAR(
14           TERMINAL[b]
15         )
16       ) # PLUS
17     ) # PAREN,
18     TERMINAL[*],
19     NUM(
20       TERMINAL[3]
21     )
22   ) # TIMES
23 ) # ASSIGN
24 ) # EXP
```

Ambigüedades y Árboles Sintácticos

Modifiquemos la precedencia de operadores utilizada en el ejemplo anterior:

```
pl@nereida:~/LEyapp/examples$ head -6 CalcSyntaxTree2.eypp | cat -n
 1 # CalcSyntaxTree.eypp
 2 %right '='
 3 %left '+'
 4 %left '-' '*' '/'
 5 %left NEG
 6 %right '^'
```

Compilamos la gramática y la ejecutamos con entrada `a=2-b*3`:

```
pl@nereida:~/LEyapp/examples$ cat prueba3.exp
a=2-b*3
pl@nereida:~/LEyapp/examples$ usecalcsyntaxtree2.pl prueba3.exp
```

```
EXP(
  ASSIGN(
    TERMINAL[a],
    TIMES(
      MINUS(
        NUM(
          TERMINAL[2]
        ),
        VAR(
          TERMINAL[b]
        )
      ) # MINUS,
      NUM(
        TERMINAL[3]
      )
    ) # TIMES
  ) # ASSIGN
) # EXP
```

Ejercicio 8.9.2. *En el programa eyapp anterior modifique las prioridades establecidas para los terminales y muestre los árboles sintácticos formados. Pruebe los siguientes experimentos.*

- *Haga el menos binario – asociativo a derechas*
- *Haga que el menos binario tenga mayor prioridad que el menos unario*
- *Introduzca en la gramática los operadores de comparación (<, >, etc.). ¿Cual es la prioridad adecuada para estos operadores? ¿Que asociatividad es correcta para los mismos? Observe los árboles formados para frases como `a = b < c * 2`. Contraste como interpreta un lenguaje típico (Java, C, Perl) una expresión como esta.*

El Método YYIssemantic

Es posible consultar o cambiar dinámicamente el estatus de un terminal usando el método `YYIssemantic` :

```
pl@nereida:~/LEyapp/examples$ sed -ne '/sub Run/, $p' CalcSyntaxTreeDynamicSemantic.eypp | cat -
1 sub Run {
2     my($self)=shift;
3
```

```

4      $input = shift;
5      $self->YYIssemantic('+', 1);
6      return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
7  }

```

Los árboles generados al ejecutar el programa contienen nodos `TERMINAL[+]`):

```

pl@nereida:~/LEyapp/examples$ usecalcsyntaxtreedynamicsemantic.pl
2+3+4

```

```

EXP(
  PLUS(
    PLUS(
      NUM(
        TERMINAL[2]
      ),
      TERMINAL[+],
      NUM(
        TERMINAL[3]
      )
    ) # PLUS,
    TERMINAL[+],
    NUM(
      TERMINAL[4]
    )
  ) # PLUS
) # EXP

```

8.10. La Maniobra de bypass

En el programa `eyapp` visto en la sección anterior generábamos una representación del árbol de análisis sintáctico utilizando la directiva `%tree`. Así para la entrada:

```

pl@nereida:~/LEyapp/examples$ cat prueba2.exp
a=(2+b)*3

```

Obteníamos el árbol:

```

pl@nereida:~/LEyapp/examples$ usecalcsyntaxtree.pl prueba2.exp | cat -n
1
2  EXP(
3    ASSIGN(
4      TERMINAL[a],
5      TIMES(
6        PAREN(
7          PLUS(
8            NUM(
9              TERMINAL[2]
10             ),
11            VAR(
12              TERMINAL[b]
13            )
14          ) # PLUS
15        ) # PAREN,

```

```

16      NUM(
17          TERMINAL[3]
18      )
19  ) # TIMES
20  ) # ASSIGN
21  ) # EXP

```

Es evidente que el árbol contiene información redundante: el nodo PAREN puede ser obviado sin que afecte a la interpretación del árbol.

La directiva `%tree` es un caso particular de la directiva `%default action`. De hecho usarla es equivalente a escribir:

```
%default action { goto &Parse::Eyapp::Driver::YYBuildAST }
```

La subrutina `Parse::Eyapp::Driver::YYBuildAST` se encarga de la construcción del nodo correspondiente a la regla de producción actual. De hecho el método retorna una referencia al objeto nodo recién construido.

Puesto que la directiva `%tree` es un caso particular de la directiva `%default action` la ubicación de acciones semánticas explícitas para una regla puede ser usada para alterar la forma del árbol de análisis. Un ejemplo de uso lo motiva el comentario anterior sobre los nodos PAREN. Sería deseable que tales nodos no formaran parte del árbol.

Observe la nueva versión de la gramática en la sección anterior (página 482). El programa `eyapp` trabaja bajo la directiva `%tree`. Se ha introducido una acción explícita en la línea 21.

```

pl@nereida:~/LEyapp/examples$ sed -ne '/^exp:/,/~/p' CalcSyntaxTree4.eyp | cat -n
 1  exp:
 2      %name NUM
 3      NUM
 4      | %name VAR
 5      VAR
 6      | %name ASSIGN
 7      VAR '=' exp
 8      | %name PLUS
 9      exp '+' exp
10      | %name MINUS
11      exp '-' exp
12      | %name TIMES
13      exp '*' exp
14      | %name DIV
15      exp '/' exp
16      | %name UMINUS
17      '-' exp %prec NEG
18      | %name EXP
19      exp '^' exp
20      | %name PAREN
21      '(' exp ')' { $_[2] }
22  ;

```

Como consecuencia de la acción de la línea 21 la referencia al nodo hijo asociado con `exp` es retornado y el nodo PAREN desaparece del árbol:

```
pl@nereida:~/LEyapp/examples$ usecalcsyntaxtree4.pl prueba2.exp
```

```

EXP(
  ASSIGN(

```



```

    TERMINAL[a],
    TIMES(
        PLUS(
            NUM(
                TERMINAL[2]
            ),
            VAR(
                TERMINAL[b]
            )
        ) # PLUS,
        NUM(
            TERMINAL[3]
        )
    ) # TIMES
) # ASSIGN
) # EXP

```

A esta maniobra, consistente en obviar un nodo sustituyéndolo por el único hijo que realmente importa, la denominaremos `bypass` de un nodo.

8.11. Salvando la Información en los Terminales Sintácticos

La directiva `%tree`

La directiva `%tree` permite que `Parse::Eyapp` genere al árbol sintáctico. Haciendo uso de las directivas adecuadas podemos controlar la forma del árbol. La directiva `%tree` es una alternativa a la directiva `%metatree` y es incompatible con esta última. El objetivo es lograr una representación adecuada del árbol sintáctico y dejar para fases posteriores la decoración del mismo.

```
nereida:~/src/perl/YappWithDefaultAction/examples> cat -n Rule6.yp
```

```

1  %{
2  use Data::Dumper;
3  %}
4  %right  '='
5  %left   '-' '+'
6  %left   '*' '/'
7  %left   NEG
8  %tree
9
10 %%
11 line: exp { $_[1] }
12 ;
13
14 exp:      %name NUM
15          NUM
16          | %name VAR
17          VAR
18          | %name ASSIGN
19          VAR '=' exp
20          | %name PLUS
21          exp '+' exp
22          | %name MINUS
23          exp '-' exp
24          | %name TIMES
25          exp '*' exp

```

```

26         | %name DIV
27         exp '/' exp
28         | %name UMINUS
29         '-' exp %prec NEG
30         | '(' exp ')' { $_[2] } /* Let us simplify a bit the tree */
31 ;
32
33 %%
34
35 sub _Error {
36     ..
37 }
38
39 sub _Lexer {
40     ..
41 }
42
43 sub Run {
44     my($self)=shift;
45     $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error,
46                   #yyprefix => "Rule6::",
47                   #yydebug =>0xFF
48                   );
49 }

```

Compilación Separada

Para compilar un programa separado `Parse::Eyapp` usamos el guión `eyapp` :

```

nereida:~/src/perl/YappWithDefaultAction/examples> eyapp Rule6
nereida:~/src/perl/YappWithDefaultAction/examples> ls -ltr | tail -1
-rw-rw----  1 pl users   5475 2006-11-06 13:53 Rule6.pm

```

Terminales sintácticos y Terminales Semánticos

A diferencia de lo que ocurre cuando se usa la directiva `%metatree` la construcción del árbol solicitado mediante la directiva `%tree` implícitamente considera token sintácticos aquellos terminales que aparecen en definidos en el programa `eyapp` mediante el uso de apóstrofes. Los token sintácticos no forman parte del árbol construido. Así pues -en el ejemplo que nos ocupa - los terminales `'=`, `'-`, `'+'`, `'*'` y `'/'` serán -por defecto - eliminados del árbol sintáctico. Por ejemplo, para la entrada `a=b*32` el siguiente árbol es construido:

```

nereida:~/src/perl/YappWithDefaultAction/examples> useruleandshift.pl
a=b*32
$VAR1 = bless( { 'children' => [
    bless( { 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL' ),
    bless( { 'children' => [
        bless( { 'children' => [
            bless( { 'children' => [], 'attr' => 'b', 'token' => 'VAR' }, 'TERMINAL' )
        ],
        }, 'VAR' ),
        bless( { 'children' => [
            bless( { 'children' => [], 'attr' => '32', 'token' => 'NUM' }, 'TERMINAL' )
        ],
        }, 'NUM' )
    ],
    }, 'VAR' )

```

```

    }, 'TIMES' )
]
}, 'ASSIGN' );

```

Esta conducta puede cambiarse usando la directiva `%semantic token` la cual declara una lista de terminales como semánticos. En tal caso dichos terminales formarán parte del árbol construido. Si por el contrario lo que queremos es cambiar el estatus de un terminal - por ejemplo `NUM` o `ID` - a sintáctico usaremos la directiva `%syntactic token`.

Salvando la Información de los `%syntactic token`

Si bien la funcionalidad de un `%syntactic token` es determinar la forma del AST, en ocasiones la diferencia entre terminal sintáctico y semántico es difusa. Un terminal fundamentalmente sintáctico pueden acarrear alguna información útil para las restantes fases de procesado.

Por ejemplo, el atributo número de línea asociado con el terminal sintáctico `'+'` puede ser útil posteriormente: Si en la fase de comprobación de tipos se observa un error de tipos incompatibles con el operador, disponer de la línea asociada con `'+'` nos permitirá emitir mensajes de error mas precisos.

En `Parse::Eyapp` el programador puede proveer a la clase `TERMINAL` de un método

```

TERMINAL::save_attributes

```

el cuál será ejecutado durante la construcción del AST cada vez que un `syntactic token` es eliminado. El método recibe como argumento - además de la referencia al nodo terminal - una referencia al nodo padre del terminal. Sigue un ejemplo:

```

sub TERMINAL::save_attributes {
    # $_[0] is a syntactic terminal
    # $_[1] is the father.
    push @{$_[1]->{lines}}, $_[0]->[1]; # save the line!
}

```

La tabla 8.1 muestra como el nodo `PLUS` recoge la información del número de línea del terminal `+`.

El árbol fué generado usando el método `str`. La información sobre los atributos se hizo mediante la siguiente subrutina:

```

sub generic_info {
    my $info;

    $info = "";
    $info .= $_[0]->type_info if $_[0]->can('type_info');

    my $sep = $info?" ":"";

    $info .= $sep.$_[0]->{line} if (defined $_[0]->{line});
    local $" = ',';
    $info .= "$sep@{$_[0]->{lines}}" if (defined $_[0]->{lines});
    return $info;
}

*PLUS::info = =>generic_info;

```

8.12. Práctica: Análisis Sintáctico

Realize durante la hora de prácticas un analizador sintáctico para el lenguaje especificado por el profesor.

Cuadro 8.1: El nodo PLUS recoge el número de línea del terminal+

Programa y Footnotes	AST generado por \$t->str
<pre> int f(int a, int b) { return a+b; } -----footnotes----- 0) Symbol Table: \$VAR1 = { 'f' => { 'type' => 'F(X_2(INT,INT),INT)', 'line' => 1 } }; ----- 1) etc.</pre>	<pre> PROGRAM^{0}(FUNCTION[f]^{1}(RETURN[2,2](PLUS[INT:2](VAR[INT](TERMINAL[a:2]), VAR[INT](TERMINAL[b:2])) # PLUS) # RETURN) # FUNCTION) # PROGRAM</pre>

8.13. Podando el Arbol

Acciones Explícitas Bajo la Directiva % tree

Como se ha mencionado anteriormente, la directiva %tree es es equivalente a escribir:

```
%default action { goto &Parse::Eyapp::Driver::YYBuildAST }
```

Donde el método `Parse::Eyapp::Driver::YYBuildAST` se encarga de retornar la referencia al objeto nodo recién construido.

La ubicación de acciones semánticas explícitas para una regla altera la forma del árbol de análisis. *El método `Parse::Eyapp::Driver::YYBuildAST` no inserta en la lista de hijos de un nodo el atributo asociado con dicho símbolo a menos que este sea una referencia. Esto conlleva que una forma de eliminar un subárbol es estableciendo una acción semántica explícita que no retorne un nodo.*

Las Declaraciones no se Incluyen en el Arbol de Análisis

El siguiente ejemplo ilustra como utilizar esta propiedad en el análisis de un lenguaje mínimo en el que los programas (P) son secuencias de declaraciones (DS) seguidas de sentencias (SS).

En un compilador típico, las declaraciones influyen en el análisis semántico (por ejemplo, enviando mensajes de error si el objeto ha sido repetidamente declarado o si es usado sin haber sido previamente declarado, etc.) y la información proveída por las mismas suele guardarse en una tabla de símbolos (declarada en la línea 5 en el siguiente código) para su utilización durante las fases posteriores.

```

pl@nereida:~/LEyapp/examples$ sed -ne '1,58p' RetUndef.eypp | cat -n
1  /* File RetUndef.eypp */
2
3  %tree
4  %{
5      use Data::Dumper;
6      my %s; # symbol table
7  %}
8
9  %%
```

```

10 P:
11     %name PROG
12     $DS $SS
13     {
14         print Dumper($DS);
15         $SS->{symboltable} = \%s;
16         return $SS;
17     }
18 ;
19
20 SS:
21     %name ST
22     S
23 | %name STS
24     SS S
25 ;
26
27 DS:
28     %name DEC
29     D
30     {}
31 | %name DECS
32     DS D
33     {}
34 ;
35
36 D : %name INTDEC
37     'int' $ID
38     {
39         die "Error: $ID declared twice\n" if exists($s{$ID});
40         $s{$ID} = 'int';
41     }
42 | %name STRDEC
43     'string' $ID
44     {
45         die "Error: $ID declared twice\n" if exists($s{$ID});
46         $s{$ID} = 'string';
47     }
48 ;
49
50 S: %name EXP
51     $ID
52     {
53         die "Error: $ID not declared\n" unless exists($s{$ID});
54         goto &Parse::Eyapp::Driver::YYBuildAST; # build the node
55     }
56 ;
57
58 %%

```

La Tabla de Símbolos

Las acciones explícitas en las líneas 37-41 y 44-47 no retornan referencias y eso significa que el subárbol de D no formará parte del árbol de análisis. En vez de ello las acciones tienen por efecto

almacenar el tipo en la entrada de la tabla de símbolos asociada con la variable.

Dando Nombre a los Atributos

Cuando en la parte derecha de una regla de producción un símbolo va precedido de un dolar como ocurre en el ejemplo con ID:

```
36 D : %name INTDEC
37     'int' $ID
38     {
39         die "Error: $ID declared twice\n" if exists($s{$ID});
40         $s{$ID} = 'int';
41     }
```

le estamos indicando a eyapp que construya una copia en una variable léxica \$ID del atributo asociado con ese símbolo. Así el fragmento anterior es equivalente a escribir:

```
D : %name INTDEC
    'int' ID
    {
        my $ID = $_[2];
        die "Error: $ID declared twice\n" if exists($s{$ID});
        $s{$ID} = 'int';
    }
```

Nótese que la tabla de símbolos es una variable léxica. Para evitar su pérdida cuando termine la fase de análisis sintáctico se guarda como un atributo del nodo programa:

```
10 P:
11     %name PROG
12     $DS $SS
13     {
14         print Dumper($DS);
15         $SS->{symboltable} = \%s;
16         return $SS;
17     }
18 ;
```

Si no se desea que el nombre del atributo coincida con el nombre de la variable se utiliza la notación *punto*. Por ejemplo:

```
exp : exp.left '+' exp.right { $left + $right }
;
```

es equivalente a:

```
exp : exp.left '+' exp.right
    {
        my $left = $_[1];
        my $right = $_[2];
        $left + $right
    }
;
```

Insertando Tareas Adicionales a la Construcción de un Nodo

Si queremos ejecutar ciertas tareas antes de la construcción de un nodo tendremos que insertar una acción semántica para las mismas. Para evitar la pérdida del subárbol deberemos llamar explícitamente a `Parse::Eyapp::Driver::YYBuildAST`. Este es el caso cuando se visitan los nodos sentencia (S):

```
50 S: %name EXP
51     $ID
52     {
53         die "Error: $ID not declared\n" unless exists($s{$ID});
54         goto &Parse::Eyapp::Driver::YYBuildAST; # build the node
55     }
56 ;
```

El Programa Cliente y su Ejecución

Veamos el programa cliente:

```
pl@nereida:~/LEyapp/examples$ cat -n useretundef.pl
1  #!/usr/bin/perl -w
2  use strict;
3  use Parse::Eyapp;
4  use RetUndef;
5  use Data::Dumper;
6
7  sub TERMINAL::info { $_[0]{attr} }
8
9  my $parser = RetUndef->new();
10 my $t = $parser->Run;
11 print $t->str,"\n";
12
13 $Data::Dumper::Indent = 1;
14 print Dumper($t);
```

Observe como no hay nodos de tipo D al ejecutar el programa con entrada `int a string b a b`:

```
pl@nereida:~/LEyapp/examples$ useretundef.pl
int a string b a b
$VAR1 = undef;
STS(ST(EXP(TERMINAL[a])),EXP(TERMINAL[b]))
```

El volcado de Dumper en la línea 14 de `RetUndef.eypp` nos muestra que el atributo asociado con DS es `undef`.

El volcado de Dumper en la línea 14 de `useretundef.pl` da el siguiente resultado:

```
$VAR1 = bless( {
  'symboltable' => {
    'a' => 'int',
    'b' => 'string'
  },
  'children' => [
    bless( {
      'children' => [
        bless( {
          'children' => [
            bless( {
              'children' => [],
```

```

        'attr' => 'a',
        'token' => 'ID'
    }, 'TERMINAL' )
]
}, 'EXP' )
]
}, 'ST' ),
bless( {
    'children' => [
        bless( {
            'children' => [],
            'attr' => 'b',
            'token' => 'ID'
        }, 'TERMINAL' )
    ]
}, 'EXP' )
]
}, 'STS' );

```

8.14. Nombres para los Atributos

Desventajas de la Notación Posicional

Dentro de las acciones los atributos de la parte derecha de la regla de producción $A \rightarrow X_1 \dots X_n$ se pasan como parámetros en $\$_{[1]}$, $\$_{[2]}$, etc. La notación posicional para los atributos puede ser inconveniente si el número de símbolos es grande. También puede resultar inconveniente durante la fase de desarrollo: Puede ocurrir que hemos escrito la regla $A \rightarrow X_1 X_2 X_3$ y la acción semántica (por ejemplo $\{ \$f = \$_{[1]} * \$_{[3]} + \$_{[2]}; \$_{[2]} \}$ y - en determinado momento - nos damos cuenta que debemos cambiar la regla añadiendo un nuevo símbolo $A \rightarrow X_1 Y X_2 X_3$ o suprimiendo alguno que ya existía. O quizá queremos insertar una acción en algún lugar intermedio de la regla. O quizá reordenamos los símbolos en la parte derecha. En todos estos caso nos encontramos en la situación de que debemos reenumerar todos los argumentos dentro de la acción semántica. En el ejemplo anterior tendríamos que reescribir la acción como:

```
{ $f = $_[1]*$_[4]+$_[3]; $_[3] }
```

La notación Dolar

Para tratar con esta situación eyapp provee mecanismos para darle nombres a los atributos semánticos asociados con los símbolos. Uno de esos mecanismos es prefijar el símbolo con un dolar. Ello indica que el nombre del símbolo puede ser usado dentro de la acción como nombre del atributo asociado. Por ejemplo, el código en la línea 22 imprime el atributo asociado con la variable sintáctica **expr**, que en este caso es su valor numérico.

La Notación Punto

La otra notación usada para darle nombres a los atributos consiste en concatenar el símbolo en la parte derecha de la regla de un punto seguido del nombre del atributo (el código completo figura en el apéndice en la página ??):

```

26  exp:
27      NUM
28      | $VAR                { $s{$VAR} }
29      | VAR.x '=' exp.y     { $s{$x} = $y }
30      | exp.x '+' exp.y     { $x + $y }
31      | exp.x '-' exp.y     { $x - $y }

```



```

32 | exp.x '*' exp.y      { $x * $y }
33 | exp.x '^' exp.y     { $x ** $y }
34 | exp.x '/' exp.y
35 | {
36 |     my $parser = shift;
37 |
38 |     $y and return($x / $y);
39 |     $parser->YYData->{ERRMSG}
40 |     = "Illegal division by zero.\n";
41 |     $parser->YYError;
42 |     undef
43 | }
44 | '-' $exp %prec NEG
45 | { -$exp }
46 | '(' $exp ')'
47 | { $exp }
48 ;

```

La Cabecera y el Arranque

El atributo asociado con **start** (línea 12) es una referencia a un par. El segundo componente del par es una referencia a la tabla de símbolos. El primero es una referencia a la lista de valores resultantes de la evaluación de las diferentes expresiones.

```

pl@nereida:~/LEyapp/examples$ cat -n CalcSimple.eyp
 1 # CalcSimple.eyp
 2 %right '='
 3 %left '-','+'
 4 %left '*','/'
 5 %left NEG
 6 %right '^'
 7 %{
 8 my %s;
 9 %}
10
11 %%
12 start: input { [ $_[1], \%s] }
13 ;
14
15 input:
16     /* empty */      { [] }
17 | input line { push(@{$_[1]},$_[2]) if defined($_[2]); $_[1] }
18 ;

```

La línea 17 indica que el atributo asociado con la variable sintáctica **input** es una referencia a una pila y que el atributo asociado con la variable sintáctica **line** debe empujarse en la pila.

De hecho, el atributo asociado con **line** es el valor de la expresión evaluada en esa línea. Así pues el atributo retornado por **input** es una referencia a una lista conteniendo los valores de las expresiones evaluadas. Observe que expresiones erróneas no aparecerán en la lista.

```

20 line:
21     '\n'          { undef }
22 | $exp '\n'      { print "$exp\n"; $exp }
23 | error '\n'    { $_[0]->YYError; undef }
24 ;

```

La Cola

Veamos el código del analizador léxico:

```
50  %%
51
52  sub _Error {
53  .. .....
62  }
63
64  my $input;
65
66  sub _Lexer {
67    my($parser)=shift;
68
69    # topicalize $input
70    for ($input) {
71      s/^[ \t]//;      # skip whites
72      return('',undef) unless $_;
73      return('NUM',$1) if s/^([0-9]+(?:\.[0-9]+)?)$/;
74      return('VAR',$1) if s/^[A-Za-z][A-Za-z0-9_]*$/;
75      return($1,$1)    if s/^(.)/s;
76    }
77
78    return('',undef);
79  }
80
81  sub Run {
82    my($self)=shift;
83
84    $input = shift;
85    return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error,
86                          #yydebug => 0xF
87                          );
88  }
```

Recuperación de Errores

Las entradas pueden contener errores. El lenguaje *eyapp* proporciona un *token* especial, **error**, que puede ser utilizado en el programa fuente para extender el analizador con producciones de error que lo dotan de cierta capacidad para recuperarse de una entrada errónea y poder continuar analizando el resto de la entrada.

```
20  line:
21    '\n'          { undef }
22    | $exp '\n'    { print "$exp\n"; $exp }
23    | error '\n'   { $_[0]->YYErrork; undef }
24  ;
```

Cuando se produce un error en el análisis, *eyapp* emite un mensaje de error (Algo como "**syntax error**") y produce "mágicamente" el terminal especial denominado **error**. A partir de ahí permanecerá silencioso, consumiendo terminales hasta encontrar una regla que sea capaz de consumir el terminal **error**. La idea general es que, a través de la regla de recuperación de errores de la línea 23 se indica que cuando se produzca un error el analizador debe descartar todos los *tokens* hasta llegar a un retorno de carro.

Además, mediante la llamada al método `$_[0]->YYError` el programador anuncia que, si se alcanza este punto, la recuperación puede considerarse "completa" y que `eyapp` puede emitir a partir de ese momento mensajes de error con la seguridad de que no son consecuencia de un comportamiento inestable provocado por el primer error.

La Acción por Defecto

Observemos la regla:

```
26  exp:
27      NUM
```

La acción por defecto es retornar `$_[1]`. Por tanto, en el caso de la regla de la línea 27 el valor retornado es el asociado a NUM.

Manejo de los símbolos

La calculadora usa un hash léxico `%s` como tabla de símbolos. Cuando dentro de una expresión encontramos una alusión a una variable retornamos el valor asociado `$s{$_[1]}` que ha sido guardado en la correspondiente entrada de la tabla de símbolos:

```
.    ...
7    %{
8    my %s;
9    %}
..   ...
11   %%
..   ...
26   exp:
27       NUM
28       | $VAR                { $s{$VAR} }
29       | VAR.x '=' exp.y     { $s{$x} = $y }
30       | exp.x '+' exp.y     { $x + $y }
..   ...
```

El Atributo YYData

En la regla de la división comprobamos que el divisor es distinto de cero.

```
34     | exp.x '/' exp.y
35     {
36         my $parser = shift;
37
38         $y and return($x / $y);
39         $parser->YYData->{ERRMSG}
40         = "Illegal division by zero.\n";
41         $parser->YYError;
42         undef
43     }
```

El método `YYData` provee acceso a un hash que se maneja como una zona de datos para el programa cliente. En el ejemplo usamos una entrada `ERRMSG` para alojar el mensaje de error.

Este mensaje es aprovechado por la subrutina de tratamiento de errores:

```
52  sub _Error {
53      my $private = $_[0]->YYData;
54
55      exists $private->{ERRMSG}
```

```

56     and do {
57         print $private->{ERRMSG};
58         delete $private->{ERRMSG};
59         return;
60     };
61     print "Syntax error.\n";
62 }

```

La subrutina `_Error` es llamada por el analizador sintáctico generado por `eyapp` cada vez que ocurre un error sintáctico. Para que ello sea posible en la llamada al analizador se especifican quienes son la rutina de análisis léxico y quien es la rutina a llamar en caso de error:

```
$self->YYParse( yylex => \&_Lexer, yyerror => \&_Error );
```

El Programa Cliente

Veamos los contenidos del ejecutable `usecalcsimple.pl` el cuál utiliza el módulo generado por `eyapp`:

```

pl@nereida:~/LEyapp/examples$ cat -n usecalcsimple.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3  use CalcSimple;
 4  use Carp;
 5
 6  sub slurp_file {
 7      my $fn = shift;
 8      my $f;
 9
10      local $/ = undef;
11      if (defined($fn)) {
12          open $f, $fn
13      }
14      else {
15          $f = \*STDIN;
16      }
17      my $input = <$f>;
18      return $input;
19  }
20
21  my $parser = CalcSimple->new();
22
23  my $input = slurp_file( shift() );
24  my ($r, $s) = @{$parser->Run($input)};
25
26  print "==== Results =====\n";
27  print "$_\n" for @$r;
28  print "==== Symbol Table =====\n";
29  print "$_ = $s->{$_}\n" for sort keys %$s;

```

Ejecución

Veamos una ejecución:

```

pl@nereida:~/LEyapp/examples$ cat prueba.exp
a=2*3

```

```

b=a+1
pl@nereida:~/LEyapp/examples$ usecalcsimple.pl prueba.exp
6
7
===== Results =====
6
7
===== Symbol Table =====
a = 6
b = 7

```

La primera aparición de la cadena "6\n7" es debida a la acción `print` en el cuerpo de la gramática. La segunda es el resultado de la línea 27 en el programa cliente.

8.15. Bypass Automático

La directiva `%tree` admite la opción `bypass`. Esta opción ayuda en la fase de construcción del AST automatizando *operaciones de bypass*. Entendemos por *operación de bypass* a la operación de retornar el nodo hijo al nodo padre del actual (abuelo del hijo), obviando así la presencia del nodo actual en el AST final.

La opción de `bypass` modifica la forma en la que Eyapp construye el árbol: Si se usa la opción `bypass`, Eyapp hará automáticamente una operación de bypass a todo nodo que resulte con un sólo hijo en tiempo de construcción el árbol. Además Eyapp cambiará la clase del nodo hijo a la clase del nodo siendo puenteado si a este se le dió un nombre explícitamente a través de la directiva `%name`.

Hay dos razones principales por las que un nodo pueda resultar - en tiempo de construcción del árbol - con un sólo hijo:

- La mas obvia es que el nodo tuviera ya un sólo hijo por la forma de la gramática: el nodo en cuestión se corresponde con la parte izquierda de una *regla de producción unaria* de la forma $E \rightarrow NUM$ o $E \rightarrow VAR$.

```

25 exp:
26     %name NUM
27     NUM
28     | %name VAR
29     VAR

```

En este caso un árbol de la forma `PLUS(NUM(TERMINAL[4]), VAR(TERMINAL[a]))` se verá transformado en `PLUS(NUM[4], VAR[a])`. El primer nodo `TERMINAL` será rebautizado (re-bendecido) en la clase `NUM`, y el segundo en la clase `VAR`.

Observe que este proceso de renombre ocurre sólo si el nodo eliminado tiene un nombre explícito dado con la directiva `%name`. Así en la regla unaria:

```

pl@nereida:~/LEyapp/examples$ sed -ne '/^line:/,/~/;$/p' TreeBypass.eyp | cat -n
1 line:
2     exp '\n'
3 ;

```

ocurre un bypass automático pero el nodo hijo de `exp` conserva su nombre. Por ejemplo, un árbol como `exp_10(PLUS(NUM(TERMINAL), NUM(TERMINAL)))` se transforma en `PLUS(NUM, NUM)`. Mientras que los nodos `TERMINAL` fueron renombrados como nodos `NUM` el nodo `PLUS` no lo fué.

- La otra razón para que ocurra un bypass automático es que - aunque la parte derecha de la regla de producción tuviera mas de un símbolo - sus restantes hijos sean podados por ser terminales sintácticos. Esto ocurre, por ejemplo, con los paréntesis en la regla $E \rightarrow (E)$.

Por ejemplo, un árbol como:

```
PAREN(PLUS(NUM( TERMINAL[2]),NUM(TERMINAL[3])))
```

se convierte en:

```
PLUS(NUM[2],NUM[3])
```

Un Ejemplo con bypass Automático

```
pl@nereida:~/LEyapp/examples$ cat -n TreeBypass.eyp
```

```

1  # TreeBypass.eyp
2  %right  '='
3  %left   '-' '+'
4  %left   '*' '/'
5  %left   NEG
6  %right  '^'
7
8  %tree bypass
9
10 %{\
11 sub NUM::info {
12     my $self = shift;
13
14     $self->{attr};
15 }
16
17 *VAR::info = \&NUM::info;
18 %}
19 %%
```

Dado que los nodos `TERMINAL` serán rebautizados como nodos `VAR` y `NUM`, dotamos a dichos nodos de métodos `info` (líneas 11-17) que serán usados durante la impresión del árbol con el método `str`.

```

21 line:
22     exp '\n'
23 ;
24
25 exp:
26     %name NUM
27     NUM
28     | %name VAR
29     VAR
30     | %name ASSIGN
31     var '=' exp
32     | %name PLUS
33     exp '+' exp
34     | %name MINUS
35     exp '-' exp
36     | %name TIMES
37     exp '*' exp
```

```

38 | %name DIV
39   exp '/' exp
40 | %name UMINUS
41   '-' exp %prec NEG
42 | %name EXPON
43   exp '^' exp
44 | '(' exp ')'
45 ;
46
47 var:
48   %name VAR
49   VAR
50 ;
51 %%

```

Observe las líneas 30-31 y 47-49. La gramática original tenía sólo una regla para la asignación:

```

exp:
  %name NUM
  NUM
  | %name VAR
  VAR
  | %name ASSIGN
  VAR '=' exp

```

Si se hubiera dejado en esta forma un árbol como `ASSIGN(TERMINAL[a], NUM(TERMINAL[4]))` quedaría como `ASSIGN(TERMINAL[a], NUM[4])`. La introducción de la variable sintáctica `var` en la regla de asignación y de su regla unaria (líneas 47-50) produce un bypass y da lugar al árbol `ASSIGN(VAR[a], NUM[4])`.

Rutinas de Soporte

```

53 sub _Error {
54     exists $_[0]->YYData->{ERRMSG}
55     and do {
56         print $_[0]->YYData->{ERRMSG};
57         delete $_[0]->YYData->{ERRMSG};
58         return;
59     };
60     print "Syntax error.\n";
61 }
62
63 my $input;
64
65 sub _Lexer {
66     my($parser)=shift;
67
68     # topicalize $input
69     for ($input) {
70         s/^[ \t]+//;      # skip whites
71         return('','undef) unless $_;
72         return('NUM',$1) if s/^[0-9]+(?:\.[0-9]+)?/{}/;
73         return('VAR',$1) if s/^[A-Za-z][A-Za-z0-9_]*/;/;
74         return($1,$1)    if s/^(.)//s;
75     }
76     return('','undef);

```

```

77 }
78
79 sub Run {
80     my($self)=shift;
81
82     $input = shift;
83     print $input;
84     return $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error,
85         #yydebug => 0xF
86     );
87 }

```

El Programa Cliente

```

pl@nereida:~/LEyapp/examples$ cat -n usetreebypass.pl
 1  #!/usr/bin/perl -w
 2  # usetreebypass.pl prueba2.exp
 3  use strict;
 4  use TreeBypass;
 5
 6  sub slurp_file {
 7      my $fn = shift;
 8      my $f;
 9
10      local $/ = undef;
11      if (defined($fn)) {
12          open $f, $fn or die "Can't find file $fn!\n";
13      }
14      else {
15          $f = \*STDIN;
16      }
17      my $input = <$f>;
18      return $input;
19  }
20
21  my $parser = TreeBypass->new();
22
23  my $input = slurp_file( shift() );
24  my $tree = $parser->Run($input);
25  die "There were errors\n" unless defined($tree);
26
27  $Parse::Eyapp::Node::INDENT = 2;
28  print $tree->str."\n";

```

Ejecuciones Veamos las salidas obtenidas sobre diversas entradas:

```

pl@nereida:~/LEyapp/examples$ eyapp TreeBypass.eyp
pl@nereida:~/LEyapp/examples$ usetreebypass.pl prueba2.exp
a=(2+b)*3

```

```

ASSIGN(
  VAR[a],
  TIMES(
    PLUS(

```



```

        NUM[2],
        VAR[b]
    ),
    NUM[3]
) # TIMES
) # ASSIGN
pl@nereida:~/LEyapp/examples$ usetreebypass.pl prueba3.exp
a=2-b*3

ASSIGN(
    VAR[a],
    MINUS(
        NUM[2],
        TIMES(
            VAR[b],
            NUM[3]
        )
    ) # MINUS
) # ASSIGN
pl@nereida:~/LEyapp/examples$ usetreebypass.pl prueba4.exp
4*3

TIMES(
    NUM[4],
    NUM[3]
)
pl@nereida:~/LEyapp/examples$ usetreebypass.pl prueba5.exp
2-)3*4
Syntax error.
There were errors

```

La Directiva no bypass La cláusula `bypass` suele producir un buen número de podas y reorganizaciones del árbol.

Es preciso tener especial cuidado en su uso. De hecho, el programa anterior contiene errores. Obsérvese la conducta del analizador para la entrada `-(2-3)`:

```

pl@nereida:~/LEyapp/examples$ usetreebypass.pl prueba7.exp
-(2-3)

UMINUS(
    NUM[2],
    NUM[3]
)

```

Que es una salida errónea: además de los `bypasses` en los terminales se ha producido un `bypass` adicional sobre el nodo `MINUS` en el árbol original

```
UMINUS(MINUS(NUM(TERMINAL[2],NUM(TERMINAL[3]))))
```

dando lugar al árbol:

```
UMINUS(NUM[2],NUM[3])
```

El `bypass` automático se produce en la regla del menos unario ya que el terminal `'-'` es por defecto - un terminal sintáctico:

```

25  exp:
..      .....
40      | %name UMINUS
41      '-' exp %prec NEG

```

Mediante la aplicación de la directiva `%no bypass UMINUS` a la regla (línea 16 abajo) inhibimos la aplicación del `bypass` a la misma:

```

pl@nereida:~/LEyapp/examples$ sed -ne '/^exp:/,/~/;$/p' TreeBypassNoBypass.eyp | cat -n
 1  exp:
 2      %name NUM
 3      NUM
 4      | %name VAR
 5      VAR
 6      | %name ASSIGN
 7      var '=' exp
 8      | %name PLUS
 9      exp '+' exp
10      | %name MINUS
11      exp '-' exp
12      | %name TIMES
13      exp '*' exp
14      | %name DIV
15      exp '/' exp
16      | %no bypass UMINUS
17      '-' exp %prec NEG
18      | %name EXPON
19      exp '^' exp
20      | '(' exp ')'
21      ;

```

```

pl@nereida:~/LEyapp/examples$ usetreebypassnobypass.pl prueba7.exp
-(2-3)

```

```

UMINUS(
  MINUS(
    NUM[2],
    NUM[3]
  )
) # UMINUS

```

El Método YYBypassrule

Aún mas potente que la directiva es usar el método `YYBypassrule` el cual permite modificar dinámicamente el estatus de `bypass` de una regla de producción. Vea esta nueva versión del anterior ejemplo:

```

pl@nereida:~/LEyapp/examples$ sed -ne '/%%/,/%%/p' TreeBypassDynamic.eyp | cat -n
 1  %%
 2
 3  line:
 4      exp '\n'
 5  ;
 6
 7  exp:
 8      %name NUM

```

```

9      NUM
10     | %name VAR
11     VAR
12     | %name ASSIGN
13     var '=' exp
14     | %name PLUS
15     exp '+' exp
16     | %name MINUS
17     exp '-' exp
18     | %name TIMES
19     exp '*' exp
20     | %name DIV
21     exp '/' exp
22     | %name UMINUS
23     '-' exp %prec NEG
24     {
25         $_[0]->YYBypassrule(0);
26         goto &Parse::Eyapp::Driver::YYBuildAST;
27     }
28     | %name EXPON
29     exp '^' exp
30     | '(' exp ')'
31 ;
32
33 var:
34     %name VAR
35     VAR
36 ;
37 %%

```

El analizador produce un árbol sintáctico correcto cuando aparecen menos unarios:

```

pl@nereida:~/LEyapp/examples$ usetreebypassdynamic.pl
-(2--3*5)
-(2--3*5)

```

```

UMINUS(
  MINUS(
    NUM[2],
    TIMES(
      UMINUS(
        NUM[3]
      ),
      NUM[5]
    ) # TIMES
  ) # MINUS
) # UMINUS

```

8.16. La opción alias de %tree

La opción `alias` de la directiva `%tree` da lugar a la construcción de métodos de acceso a los hijos de los nodos para los cuales se haya explicitado un nombre a través de la *notación punto* o de la *notación dolar*.

Veamos un ejemplo:

```

nereida:~/src/perl/YappWithDefaultAction/examples> cat -n ./alias.pl
 1  #!/usr/bin/perl -w
 2  use strict;
 3  use Parse::Eyapp;
 4
 5  my $grammar = q{
 6  %right  '='
 7  %left   '-' '+'
 8  %left   '*' '/'
 9  %left   NEG
10  %tree bypass alias
11
12  %%
13  line: $exp { $_[1] }
14  ;
15
16  exp:
17      %name NUM
18          $NUM
19      | %name VAR
20          $VAR
21      | %name ASSIGN
22          $VAR '=' $exp
23      | %name PLUS
24          exp.left '+' exp.right
25      | %name MINUS
26          exp.left '-' exp.right
27      | %name TIMES
28          exp.left '*' exp.right
29      | %name DIV
30          exp.left '/' exp.right
31      | %no bypass UMINUS
32          '-' $exp %prec NEG
33  | '(' exp ')' { $_[2] } /* Let us simplify a bit the tree */
34  ;
35
36  %%
37
38  sub _Error {
39      exists $_[0]->YYData->{ERRMSG}
40      and do {
41          print $_[0]->YYData->{ERRMSG};
42          delete $_[0]->YYData->{ERRMSG};
43          return;
44      };
45      print "Syntax error.\n";
46  }
47
48  sub _Lexer {
49      my($parser)=shift;
50
51          $parser->YYData->{INPUT}
52      or $parser->YYData->{INPUT} = <STDIN>

```

```

53     or return('',undef);
54
55     $parser->YYData->{INPUT} =~ s/^ \s+ //;
56
57     for ($parser->YYData->{INPUT}) {
58         s/^( [0-9]+ (?: \. [0-9]+ )? ) //
59         and return('NUM',$1);
60         s/^( [A-Za-z] [A-Za-z0-9_]* ) //
61         and return('VAR',$1);
62         s/^( . ) //s
63         and return($1,$1);
64     }
65 }
66
67 sub Run {
68     my($self)=shift;
69     $self->YYParse( yylex => \&_Lexer, yyerror => \&_Error,
70                   #yydebug =>0xFF
71                   );
72 }
73 }; # end grammar
74
75
76 Parse::Eyapp->new_grammar(
77     input=>$grammar,
78     classname=>'Alias',
79     firstline =>7,
80 );
81 my $parser = Alias->new();
82 $parser->YYData->{INPUT} = "a = -(2*3+5-1)\n";
83 my $t = $parser->Run;
84 $Parse::Eyapp::Node::INDENT=0;
!85 print $t->VAR->str."\n";           # a
86 print "*****\n";
!87 print $t->exp->exp->left->str."\n"; # 2*3+5
88 print "*****\n";
!89 print $t->exp->exp->right->str."\n"; # 1

```

Este programa produce la salida:

```

nereida:~/src/perl/YappWithDefaultAction/examples> ./alias.pl
TERMINAL
*****
PLUS(TIMES(NUM,NUM),NUM)
*****
NUM

```

El efecto de la opción `alias` en una regla como:

```

27         | %name TIMES
28         exp.left '*' exp.right

```

es crear dos métodos `left` y `right` en la clase `TIMES`. El método `left` permite acceder al hijo izquierdo del nodo en el AST y el método `right` al hijo derecho.

El efecto de la opción `alias` en una regla como:

```
21      | %name ASSIGN
22      $VAR '=' $exp
```

es crear dos métodos en la clase `ASSIGN` denominados `VAR` y `exp` los cuales permiten acceder a los correspondientes hijos de uno nodo `ASSIGN`.

Los métodos son construidos durante la compilación de la gramática. Por tanto, las modificaciones al AST que el programa introduzca en tiempo de ejecución (por ejemplo la supresión de hijos de ciertos nodos) invalidarán el uso de estos métodos.

8.17. Diseño de Analizadores con Parse::Eyapp

A la hora de construir un analizador sintáctico tenga en cuenta las siguientes normas de buena programación:

1. Comience trabajando en el cuerpo de la gramática.
2. Olvídense al principio del analizador léxico. Su primer objetivo es tener una gramática limpia de conflictos y que reconozca el lenguaje dado.
3. Sustituya las repeticiones BNF por listas usando los operadores **eyapp** `+`, `*` y sus variantes con separadores. Si una variable describe una lista de **cosas** pongale un adjetivo adecuado como **cosaslist**. Ponga nombres significativos a las variables y terminales. No los llame **d1**, **d2**, etc.
4. Si tiene un elemento opcional en la BNF, por ejemplo, en la regla:

functiondefinition \rightarrow [basictype] functionheader functionbody
use el operador ?.

5. Cada par de reglas que introduzca vuelva a recompilar con **eyapp** la gramática para ver si se introducido ambigüedad. Cuando estoy editando la gramática suelo escribir a menudo la orden
- ```
:!eyapp %
```
- para recompilar:

[illegible]

Esto llama a **eyapp** con el fichero bajo edición. Si hay errores o conflictos (esto es, hemos introducido ambigüedad) los detectaremos enseguida. Procure detectar la aparición de un conflicto lo antes posible. Observe el sangrado del ejemplo. Es el que le recomiendo.

6. Cuando esté en el proceso de construcción de la gramática y aún le queden por rellenar variables sintácticas, declárelas como terminales mediante **%token**. De esta manera evitará las quejas de **eyapp**.

## 7. Resolución de Ambigüedades y Conflictos

Las operaciones de asignación tienen la prioridad mas baja, seguidas de las lógicas, los test de igualdad, los de comparación, a continuación las aditivas, multiplicativas y por último las operaciones de tipo **unary** y **primary**. Expresa la asociatividad natural y la prioridad especificada usando los mecanismos que **eyapp** provee al efecto: **%left**, **%right**, **%nonassoc** y **%prec**.

8. La gramática de SimpleC es ambigua, ya que para una sentencia como

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$$

existen dos árboles posibles: uno que asocia el “else” con el primer “if” y otra que lo asocia con el segundo. Los dos árboles corresponden a las dos posibles parentizaciones:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1 \text{ else } S_2)$$

Esta es la regla de prioridad usada en la mayor parte de los lenguajes: un “else” casa con el “if” mas cercano. La otra posible parentización es:

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1) \text{ else } S_2$$

*La conducta por defecto de **eyapp** es parentizar a derechas.* El generador **eyapp** nos informará del conflicto pero si no se le indica como resolverlo parentizará a derechas. Resuelva este conflicto.

9. *¿Que clase de árbol debe producir el analizador?* La respuesta es que sea lo mas abstracto posible. Debe

- Contener toda la información necesaria para el manejo eficiente de las fases subsiguientes: Análisis de ámbito, Comprobación de tipos, Optimización independiente de la máquina, etc.
- Ser uniforme
- Legible (human-friendly)
- No contener nodos que no portan información.

El siguiente ejemplo muestra una versión aceptable de árbol abstracto. Cuando se le proporciona el programa de entrada:

```
nereida:~/doc/casiano/PLB00K/PLB00K/code> cat -n prueba5.c
1 int f(int a)
2 {
3 if (a>0)
4 a = f(a-1);
5 }
```

El siguiente árbol ha sido producido por un analizador usando la directiva **%tree** y añadiendo las correspondientes acciones de **bypass**. Puede considerarse un ejemplo aceptable de AST:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> eyapp Simple2 ;\
 usesimple2.pl prueba5.c
```

```
PROGRAM(
 TYPEDFUNC(
 INT(TERMINAL[INT:1]),
 FUNCTION(
 TERMINAL[f:1],
 PARAMS(
 PARAM(
 INT(TERMINAL[INT:1]),
 TERMINAL[a:1],
 ARRAYSPEC
)
),
 BLOCK(
 DECLARATIONS,
 STATEMENTS(
 IF(
 GT(
 VAR(TERMINAL[a:3]),
 INUM(TERMINAL[0:3])
),
 ASSIGN(
 VAR(TERMINAL[a:4]),
 FUNCTIONCALL(
 TERMINAL[f:4],
 ARGLIST(
 MINUS(
 VAR(TERMINAL[a:4]),
 INUM(TERMINAL[1:4])
)
) # ARGLIST
) # FUNCTIONCALL
) # ASSIGN
) # IF
) # STATEMENTS
) # BLOCK
) # FUNCTION
) # TYPEDFUNC
) # PROGRAM
```

Es deseable darle una estructura uniforme al árbol. Por ejemplo, como consecuencia de que la gramática admite funciones con declaración implícita del tipo retornado cuando este es entero

```
1 definition:
2 funcDef { $_[1]->type("INTFUNC"); $_[1] }
3 | %name TYPEDFUNC
4 basictype funcDef
5 | declaration { $_[1] }
6 ;
```

se producen dos tipos de árboles. Es conveniente convertir las definiciones de función con declaración implícita en el mismo árbol que se obtiene con declaración explícita.



## 8.18. Práctica: Construcción del Arbol para el Lenguaje Simple

Reescriba la gramática del lenguaje Simple introducido en la práctica 8.7 para producir un árbol abstracto fácil de manejar. La gramática original produce árboles excesivamente profundos y profusos. Reescriba la gramática usando precedencia de operadores para deshacer las ambigüedades. Optimice la forma de los árboles usando bypass donde sea conveniente.

Las declaraciones no formaran parte del árbol. Utilice una tabla de símbolos para guardar la información aportada en las declaraciones sobre los identificadores. Compruebe que ninguna variable es declarada dos veces y que no se usan variables sin declarar.

## 8.19. Práctica: Ampliación del Lenguaje Simple

Amplíe la versión del lenguaje Simple desarrollada en la práctica 8.18 para que incluya sentencias `if`, `if ... else ...`, bucles `while` y `do ... while`.

Añada el tipo array. Deberían aceptarse declaraciones como esta:

```
[10][20]int a,b;
[5]string c;
```

y sentencias como esta:

```
a[4][7] = b[2][1]*2;
c[1] = "hola";
```

Diseñe una representación adecuada para almacenar los tipos array en la tabla de símbolos.

## 8.20. Agrupamiento y Operadores de Listas

Los operadores de listas `*`, `+` y `?` pueden usarse en las partes derechas de las reglas de producción para indicar repeticiones<sup>1</sup>.

### El Operador Cierre Positivo

La gramática:

```
pl@nereida:~/LEyapp/examples$ head -12 List3.y | cat -n
1 # List3.y
2 %semantic token 'c'
3 %{
4 use Data::Dumper;
5 %}
6 %%
7 S: 'c'+ 'd'+
8 {
9 print Dumper($_[1]);
10 print Dumper($_[2]);
11 }
12 ;
```

Es equivalente a:

```
pl@nereida:~/LEyapp/examples$ eyapp -v List3.y | head -9 List3.output
Rules:

```

---

<sup>1</sup>La descripción que se hace en esta versión requiere una versión de `eyapp` igual o posterior a 1.087

```

0: $start -> S $end
1: PLUS-1 -> PLUS-1 'c'
2: PLUS-1 -> 'c'
3: PLUS-2 -> PLUS-2 'd'
4: PLUS-2 -> 'd'
5: S -> PLUS-1 PLUS-2

```

La acción semántica asociada con un operador de lista + retorna una lista con los atributos de los factores a los que afecta el +:

```

pl@nereida:~/LEyapp/examples$ use_list3.pl
ccdd
$VAR1 = ['c', 'c'];
$VAR1 = ['d', 'd'];

```

Si se activado una de las directivas de creación de árbol (%tree o %metatree) o bien se ha hecho uso explícito del método YYBuildingTree o se pasa la opción yybuildingtree de YYParse la semántica de la acción asociada cambia. En tal caso la acción semántica asociada con un operador de lista + es crear un nodo con la etiqueta

\_PLUS\_LIST\_#number

cuyos hijos son los elementos de la lista. El número es el número de orden de la regla tal y como aparece en el fichero .output. Como ocurre cuando se trabaja bajo la directiva %tree, es necesario que los atributos de los símbolos sean terminales semánticos o referencias para que se incluyan en la lista.

Al ejecutar el programa anterior bajo la directiva %tree se obtiene la salida:

```

pl@nereida:~/LEyapp/examples$ head -3 List3.y; eyapp List3.y
List3.y
%semantic token 'c'
%tree

pl@nereida:~/LEyapp/examples$ use_list3.pl
ccdd
$VAR1 = bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'c', 'token' => 'c' }, 'TERMINAL'),
 bless({ 'children' => [], 'attr' => 'c', 'token' => 'c' }, 'TERMINAL')
],
 }, '_PLUS_LIST_1');
$VAR1 = bless({ 'children' => [] }, '_PLUS_LIST_2');

```

El nodo asociado con la lista de ds aparece vacío por que el terminal 'd' no fué declarado semántico.

## Desaparición de Nodos en Listas

Cuando se trabaja con listas bajo la directiva %tree la acción por defecto es el "aplanamiento" de los nodos a los que se aplica el operador + en una sólo lista.

En el ejemplo anterior los nodos 'd' no aparecen pues 'd' es un token sintáctico. Sin embargo, puede que no sea suficiente declarar 'd' como semántico.

Cuando se construye el árbol, el algoritmo de construcción de nodos asociados con las listas omite cualquier atributo que no sea una referencia. Por tanto, es necesario garantizar que el atributo asociado con el símbolo sea una referencia (o un token semántico) para asegurar su presencia en la lista de hijos del nodo.

```

pl@nereida:~/LEyapp/examples$ head -19 ListWithRefs1.eyy | cat -n
1 # ListWithRefs.eyy

```



Aunque 'd' es declarado semántico la acción por defecto asociada con la regla D: 'd' en la línea 16 retornará \$\_[1] (esto es, el escalar 'd'). Dado que no se trata de una referencia no es insertado en la lista de hijos del nodo \_PLUS\_LIST.

## Recuperando los Nodos Desaparecidos

La solución es hacer que cada uno de los símbolos a los que afecta el operador de lista sea una referencia.

```
pl@nereida:~/LEyapp/examples$ head -22 ListWithRefs.eypp | cat -n
1 # ListWithRefs.eypp
2 %semantic token 'c'
3 %{
4 use Data::Dumper;
5 %}
6 %%
7 S: 'c'+ D+
8 {
9 print Dumper($_[1]);
10 print $_[1]->str."\n";
11 print Dumper($_[2]);
12 print $_[2]->str."\n";
13 }
14 ;
15
16 D: 'd'
17 {
18 bless { attr => $_[1], children =>[] }, 'DES';
19 }
20 ;
21
22 %%
```

Ahora el atributo asociado con D es un nodo y aparece en la lista de hijos del nodo \_PLUS\_LIST:

```
pl@nereida:~/LEyapp/examples$ eyapp ListWithRefs.eypp; use_listwithrefs.pl
ccdd
$VAR1 = bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => 'c',
 'token' => 'c'
 }, 'TERMINAL'),
 bless({
 'children' => [],
 'attr' => 'c',
 'token' => 'c'
 }, 'TERMINAL')
],
 }, '_PLUS_LIST_1');
_PLUS_LIST_1(TERMINAL,TERMINAL)
$VAR1 = bless({
 'children' => [
 bless({
```

```

 'children' => [],
 'attr' => 'd'
 }, 'DES'),
 bless({
 'children' => [],
 'attr' => 'd'
 }, 'DES')
]
 }, '_PLUS_LIST_2');
_PLUS_LIST_2(DES,DES)

```

### Construcción de un Árbol con `Parse::Eyapp::Node->new`

La solución anterior consistente en escribir a mano el código de construcción del nodo puede ser suficiente cuando se trata de un sólo nodo. Escribir a mano el código para la construcción de un árbol con varios nodos puede ser tedioso. Peor aún: aunque el nodo construido en el ejemplo luce como los nodos `Parse::Eyapp` no es realmente un nodo `Parse::Eyapp`. Los nodos `Parse::Eyapp` siempre heredan de la clase `Parse::Eyapp::Node` y por tanto tienen acceso a los métodos definidos en esa clase. La siguiente ejecución con el depurador ilustra este punto:

```
pl@nereida:~/LEyapp/examples$ perl -wd use_listwithrefs.pl
```

```
Loading DB routines from perl5db.pl version 1.28
Editor support available.
```

```
Enter h or 'h h' for help, or 'man perldebug' for more help.
```

```
main: (use_listwithrefs.pl:4): $parser = new ListWithRefs();
DB<1> f ListWithRefs.eyp
1 2 #line 3 "ListWithRefs.eyp"
3
4: use Data::Dumper;
5
6 #line 7 "ListWithRefs.eyp"
7 #line 8 "ListWithRefs.eyp"
8
9: print Dumper($_[1]);
10: print $_[1]->str."\n";
```

Mediante el comando `f ListWithRefs.eyp` le indicamos al depurador que los subsiguientes comandos harán referencia a dicho fichero. A continuación ejecutamos el programa hasta alcanzar la acción semántica asociada con la regla `S: 'c'+ D+` (línea 9)

```
DB<2> c 9 # Continuar hasta la línea 9 de ListWithRefs.eyp
ccdd
ListWithRefs::CODE(0x84ebe5c)(ListWithRefs.eyp:9):
9: print Dumper($_[1]);
```

Estamos en condiciones ahora de observar los contenidos de los argumentos:

```
DB<3> x $_[2]->str
0 '_PLUS_LIST_2(DES,DES)'
DB<4> x $_[2]->child(0)
0 DES=HASH(0x85c4568)
 'attr' => 'd'
 'children' => ARRAY(0x85c458c)
 empty array
```

El método `str` funciona con el objeto `$_[2]` pues los nodos `_PLUS_LIST_2` heredan de la clase `Parse::Eyapp::Node`. sin embargo, cuando intentamos usarlo con un nodo `DES` obtenemos un error:

```
DB<6> x $_[2]->child(0)->str
Can't locate object method "str" via package "DES" at \
(eval 11) [/usr/share/perl/5.8/perl5db.pl:628] line 2, <STDIN> line 1.
DB<7>
```

Una solución mas robusta que la anterior es usar el constructor `Parse::Eyapp::Node->new`. El método `Parse::Eyapp::Node->new` se usa para construir un bosque de árboles sintácticos. Recibe una secuencia de términos describiendo los árboles y - opcionalmente - una subrutina. La subrutina se utiliza para iniciar los atributos de los nodos recién creados. Después de la creación del árbol la subrutina es llamada por `Parse::Eyapp::Node->new` pasándole como argumentos la lista de referencias a los nodos en el orden en el que aparecen en la secuencia de términos de izquierda a derecha. `Parse::Eyapp::Node->new` retorna una lista de referencias a los nodos recién creados, en el orden en el que aparecen en la secuencia de términos de izquierda a derecha. En un contexto escalar devuelve la referencia al primero de esos árboles. Por ejemplo:

```
pl@nereida:~/LEyapp/examples$ perl -MParse::Eyapp -MData::Dumper -wde 0
main::(-e:1): 0
DB<1> @t = Parse::Eyapp::Node->new('A(C,D) E(F)', sub { my $i = 0; $_->{n} = $i++ for @_ })
DB<2> $Data::Dumper::Indent = 0
DB<3> print Dumper($_)."\n" for @t
$VAR1 = bless({'n' => 0, 'children' => [bless({'n' => 1, 'children' => []}, 'C'),
 bless({'n' => 2, 'children' => []}, 'D')
]
 }, 'A');
$VAR1 = bless({'n' => 1, 'children' => []}, 'C');
$VAR1 = bless({'n' => 2, 'children' => []}, 'D');
$VAR1 = bless({'n' => 3, 'children' => [bless({'n' => 4, 'children' => []}, 'F')]}, 'E');
$VAR1 = bless({'n' => 4, 'children' => []}, 'F');
```

Véase el siguiente ejemplo en el cual los nodos asociados con las 'd' se han convertido en un subárbol:

```
pl@nereida:~/LEyapp/examples$ head -28 ListWithRefs2.eyp | cat -n
 1 # ListWithRefs2.eyp
 2 %semantic token 'c'
 3 %{
 4 use Data::Dumper;
 5 %}
 6 %%
 7 S: 'c'+ D+
 8 {
 9 print Dumper($_[1]);
10 print $_[1]->str."\n";
11 print Dumper($_[2]);
12 print $_[2]->str."\n";
13 }
14 ;
15
16 D: 'd'.d
17 {
18 Parse::Eyapp::Node->new(
19 'DES(TERMINAL)',
```

```

20 sub {
21 my ($DES, $TERMINAL) = @_ ;
22 $TERMINAL->{attr} = $d;
23 }
24);
25 }
26 ;
27
28 %%

```

Para saber mas sobre `Parse::Eyapp::Node->new` consulte la entrada `Parse::Eyapp::Node->new` de la documentación de `Parse::Eyapp`.

```

pl@nereida:~/LEyapp/examples$ eyapp ListWithRefs2.eyp; use_listwithrefs2.pl
ccdd
$VAR1 = bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'c', 'token' => 'c' }, 'TERMINAL'),
 bless({ 'children' => [], 'attr' => 'c', 'token' => 'c' }, 'TERMINAL')
]
}, '_PLUS_LIST_1');
_PLUS_LIST_1(TERMINAL,TERMINAL)
$VAR1 = bless({
 'children' => [
 bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'd' }, 'TERMINAL')
]
 }, 'DES'),
 bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'd' }, 'TERMINAL')
]
 }, 'DES')
]
}, '_PLUS_LIST_2');
_PLUS_LIST_2(DES(TERMINAL),DES(TERMINAL))

```

## El Operador de Cierre \*

Un operador de lista afecta al factor a su izquierda. Una lista en la parte derecha de una regla cuenta como un único símbolo.

Los operadores `*` y `+` pueden ser usados con el formato `X <* Separator>`. En ese caso lo que se describen son listas separadas por el separador `separator`.

```

pl@nereida:~/LEyapp/examples$ head -25 CsBetweenCommansAndD.eyp | cat -n
1 # CsBetweenCommansAndD.eyp
2
3 %semantic token 'c' 'd'
4
5 %{
6 sub TERMINAL::info {
7 $_[0]->attr;
8 }
9 %}

```

```

10 %tree
11 %%
12 S:
13 ('c' <* ','> 'd')*
14 {
15 print "\nNode\n";
16 print $_[1]->str."\n";
17 print "\nChild 0\n";
18 print $_[1]->child(0)->str."\n";
19 print "\nChild 1\n";
20 print $_[1]->child(1)->str."\n";
21 $_[1]
22 }
23 ;
24
25 %%

```

La regla S: ('c' <\* ','> 'd')\* tiene dos elementos en su parte derecha: la lista de cs separadas por comas y la d. La regla es equivalente a:

```

pl@nereida:~/LEyapp/examples$ eyapp -v CsBetweenCommansAndD.eyp
pl@nereida:~/LEyapp/examples$ head -11 CsBetweenCommansAndD.output | cat -n
 1 Rules:
 2 -----
 3 0: $start -> S $end
 4 1: STAR-1 -> STAR-1 ',' 'c'
 5 2: STAR-1 -> 'c'
 6 3: STAR-2 -> STAR-1
 7 4: STAR-2 -> /* empty */
 8 5: PAREN-3 -> STAR-2 'd'
 9 6: STAR-4 -> STAR-4 PAREN-3
10 7: STAR-4 -> /* empty */
11 8: S -> STAR-4

```

La acción semántica asociada con un operador de lista \* es retornar una referencia a una lista con los atributos de los elementos a los que se aplica.

Si se trabaja - como en el ejemplo - bajo una directiva de creación de árbol retorna un nodo con la etiqueta `_STAR_LIST_#number` cuyos hijos son los elementos de la lista. El número es el número de orden de la regla tal y como aparece en el fichero `.output`. Es necesario que los elementos sean terminales o referencias para que se incluyan en la lista. Observe como el nodo PAREN-3 ha sido eliminado del árbol. Los nodos paréntesis son -en general - obviados:

```

pl@nereida:~/LEyapp/examples$ use_csbetweencommansandd.pl
c,c,cd

```

Node

```
_STAR_LIST_4(_STAR_LIST_1(TERMINAL[c],TERMINAL[c],TERMINAL[c]),TERMINAL[d])
```

Child 0

```
_STAR_LIST_1(TERMINAL[c],TERMINAL[c],TERMINAL[c])
```

Child 1

```
TERMINAL[d]
```

Obsérvese también que la coma ha sido eliminada.



## Poniendo Nombres a las Listas

Para poner nombre a una lista la directiva `%name` debe preceder al operador tal y como ilustra el siguiente ejemplo:

```
pl@nereida:~/LEyapp/examples$ sed -ne '1,27p' CsBetweenCommansAndDWithNames.eyp | cat -n
1 # CsBetweenCommansAndDWithNames.eyp
2
3 %semantic token 'c' 'd'
4
5 %{
6 sub TERMINAL::info {
7 $_[0]->attr;
8 }
9 %}
10 %tree
11 %%
12 Start: S
13 ;
14 S:
15 ('c' <%name Cs * ','> 'd') %name Cs_and_d *
16 {
17 print "\nNode\n";
18 print $_[1]->str."\n";
19 print "\nChild 0\n";
20 print $_[1]->child(0)->str."\n";
21 print "\nChild 1\n";
22 print $_[1]->child(1)->str."\n";
23 $_[1]
24 }
25 ;
26
27 %%
```

La ejecución muestra que los nodos de lista han sido renombrados:

```
pl@nereida:~/LEyapp/examples$ use_csbetweencommansanddwithnames.pl
c,c,c,cd
```

Node

```
Cs_and_d(Cs(TERMINAL[c],TERMINAL[c],TERMINAL[c],TERMINAL[c]),TERMINAL[d])
```

Child 0

```
Cs(TERMINAL[c],TERMINAL[c],TERMINAL[c],TERMINAL[c])
```

Child 1

```
TERMINAL[d]
```

## Opcionales

La utilización del operador `?` indica la presencia o ausencia del factor a su izquierda. La gramática:

```
pl@nereida:~/LEyapp/examples$ head -11 List5.yyp | cat -n
1 %semantic token 'c'
2 %tree
3 %%
4 S: 'c' 'c'?
```

```

5 {
6 print $_[2]->str."\n";
7 print $_[2]->child(0)->attr."\n" if $_[2]->children;
8 }
9 ;
10
11 %%

```

produce la siguiente gramática:

```

l@nereida:~/LEyapp/examples$ eyapp -v List5
pl@nereida:~/LEyapp/examples$ head -7 List5.output
Rules:

0: $start -> S $end
1: OPTIONAL-1 -> 'c'
2: OPTIONAL-1 -> /* empty */
3: S -> 'c' OPTIONAL-1

```

Cuando no se trabaja bajo directivas de creación de árbol el atributo asociado es una lista (vacía si el opcional no aparece). Bajo la directiva `%tree` el efecto es crear el nodo:

```

pl@nereida:~/LEyapp/examples$ use_list5.pl
cc
_OPTIONAL_1(TERMINAL)
c
pl@nereida:~/LEyapp/examples$ use_list5.pl
c
_OPTIONAL_1

```

## Agrupamientos

Es posible agrupar mediante paréntesis una subcadena de la parte derecha de una regla de producción. La introducción de un paréntesis implica la introducción de una variable adicional cuya única producción es la secuencia de símbolos entre paréntesis. Así la gramática:

```

pl@nereida:~/LEyapp/examples$ head -6 Parenthesis.eyp | cat -n
1 %%
2 S:
3 ('a' S) 'b' { shift; [@_] }
4 | 'c'
5 ;
6 %%

```

genera esta otra gramática:

```

pl@nereida:~/LEyapp/examples$ eyapp -v Parenthesis.eyp; head -6 Parenthesis.output
Rules:

0: $start -> S $end
1: PAREN-1 -> 'a' S
2: S -> PAREN-1 'b'
3: S -> 'c'

```

Por defecto, la regla semántica asociada con un paréntesis consiste en formar una lista con los atributos de los símbolos entre paréntesis:

```
pl@nereida:~/LEyapp/examples$ cat -n use_parenthesis.pl
1 #!/usr/bin/perl -w
2 use Parenthesis;
3 use Data::Dumper;
4
5 $Data::Dumper::Indent = 1;
6 $parser = Parenthesis->new();
7 print Dumper($parser->Run);
```

```
pl@nereida:~/LEyapp/examples$ use_parenthesis.pl
```

```
acb
```

```
$VAR1 = [
 ['a', 'c'], 'b'
];
```

```
pl@nereida:~/LEyapp/examples$ use_parenthesis.pl
```

```
aacbb
```

```
$VAR1 = [
 [
 'a',
 [['a', 'c'], 'b']
],
 'b'
];
```

Cuando se trabaja bajo una directiva de creación de árbol o se establece el atributo con el método YYBuildingtree la acción semántica es construir un nodo con un hijo por cada atributo de cada símbolo entre paréntesis. Si el atributo no es una referencia no es insertado como hijo del nodo.

Veamos un ejemplo:

```
pl@nereida:~/LEyapp/examples$ head -23 List2.y | cat -n
```

```
1 %{
2 use Data::Dumper;
3 %}
4 %semantic token 'a' 'b' 'c'
5 %tree
6 %%
7 S:
8 (%name AS 'a' S)'b'
9 {
10 print "S -> ('a' S)'b'\n";
11 print "Atributo del Primer Símbolo:\n".Dumper($_[1]);
12 print "Atributo del Segundo símbolo: $_[2]\n";
13 $_[0]->YYBuildAST(@_[1..$#_]);
14 }
15 | 'c'
16 {
17 print "S -> 'c'\n";
18 my $r = Parse::Eyapp::Node->new(qw(C(TERMINAL)), sub { $_[1]->attr('c') }) ;
19 print Dumper($r);
20 $r;
21 }
22 ;
23 %%
```

El ejemplo muestra (línea 8) como renombrar un nodo \_PAREN asociado con un paréntesis. Se escribe la directiva %name CLASSNAME después del paréntesis de apertura.

La llamada de la línea 13 al método `YYBuildAST` con argumentos los atributos de los símbolos de la parte derecha se encarga de retornar el nodo describiendo la regla de producción actual. Observe que la línea 13 puede ser reescrita como:

```
goto &Parse::Eyapp::Driver::YYBuildAST;
```

En la línea 18 se construye explícitamente un nodo para la regla usando `Parse::Eyapp::Node->new`. El manejador pasado como segundo argumento se encarga de establecer un valor para el atributo `attr` del nodo `TERMINAL` que acaba de ser creado.

Veamos una ejecución:

```
pl@nereida:~/LEyapp/examples$ use_list2.pl
aacbb
S -> 'c'
$VAR1 = bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => 'c'
 }, 'TERMINAL')
]
}, 'C');
```

La primera reducción ocurre por la regla no recursiva. La ejecución muestra el árbol construido mediante la llamada a `Parse::Eyapp::Node->new` de la línea 18.

La ejecución continúa con una reducción o anti-derivación por la regla `S -> ('a' S )'b'`. La acción de las líneas 9-14 hace que se muestre el atributo asociado con `('a' S)` o lo que es lo mismo con la variable sintáctica adicional `PAREN-1` así como el atributo de `'b'`:

```
S -> ('a' S)'b'
Atributo del Primer Símbolo:
$VAR1 = bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'a', 'token' => 'a' }, 'TERMINAL'),
 bless({ 'children' => [bless({ 'children' => [], 'attr' => 'c' }, 'TERMINAL')
]
 }, 'C')
]
}, 'AS');
Atributo del Segundo símbolo: b
```

La última reducción visible es por la regla `S -> ('a' S )'b'`:

```
S -> ('a' S)'b'
Atributo del Primer Símbolo:
$VAR1 = bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'a', 'token' => 'a' }, 'TERMINAL'),
 bless({
 'children' => [
 bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'a', 'token' => 'a' }, 'TERMINAL'),
 bless({
 'children' => [
 bless({ 'children' => [], 'attr' => 'c' }, 'TERMINAL')
]
 }
]
 }
]
 }
], 'C')
]
}, 'AS');
```

```

]
 }, 'C')
]
 }, 'AS'),
 bless({ 'children' => [], 'attr' => 'b', 'token' => 'b' }, 'TERMINAL')
]
}, 'S_2')
]
}, 'AS');
Atributo del Segundo símbolo: b

```

### Acciones Entre Paréntesis

Aunque no es una práctica recomendable es posible introducir acciones en los paréntesis como en este ejemplo:

```

pl@nereida:~/LEyapp/examples$ head -16 ListAndAction.eyp | cat -n
1 # ListAndAction.eyp
2 %{
3 my $num = 0;
4 %}
5
6 %%
7 S: 'c'
8 {
9 print "S -> c\n"
10 }
11 | ('a' {$num++; print "Seen <$num> 'a's\n"; $_[1] }) S 'b'
12 {
13 print "S -> (a) S b\n"
14 }
15 ;
16 %%

```

Al ejecutar el ejemplo con la entrada `aaacbbb` se obtiene la siguiente salida:

```

pl@nereida:~/LEyapp/examples$ use_listandaction.pl
aaacbbb
Seen <1> 'a's
Seen <2> 'a's
Seen <3> 'a's
S -> c
S -> (a) S b
S -> (a) S b
S -> (a) S b

```

## 8.21. El método `str` en Mas Detalle

### El método `str` de los Nodos

Para dar soporte al análisis de los árboles y su representación, `Parse::Eyapp` provee el método `str`. El siguiente ejemplo con el depurador muestra el uso del método `str`. El método `Parse::Eyapp::Node::str` retorna una cadena que describe como término el árbol enraizado en el nodo que se ha pasado como argumento.

## El método `info`

El método `str` cuando visita un nodo comprueba la existencia de un método `info` para la clase del nodo. Si es así el método será llamado. Obsérvese como en las líneas 5 y 6 proveemos de métodos `info` a los nodos `TERMINAL` y `FUNCTION`. La consecuencia es que en la llamada de la línea 7 el término que describe al árbol es decorado con los nombres de funciones y los atributos y números de línea de los terminales.

## Modo de uso de `str`

El método `str` ha sido concebido como una herramienta de ayuda a la depuración y verificación de los programas árbol. La metodología consiste en incorporar las pequeñas funciones `info` al programa en el que trabajamos. Por ejemplo:

```
640 sub Run {
641 my($self)=shift;
642
643 }
644
645 sub TERMINAL::info {
646 my @a = join ':', @{$_[0]->{attr}};
647 return "[@a]"
648 }
649
650 sub FUNCTION::info {
651 return "[".$_[0]->{function_name}[0]."]"
652 }
653
654 sub BLOCK::info {
655 return "[".$_[0]->{line}."]"
656 }
```

## Variables que Gobiernan la Conducta de `str`

Las variables de paquete que gobiernan la conducta de `str` y sus valores por defecto aparecen en la siguiente lista:

- `our @PREFIXES = qw(Parse::Eyapp::Node::)`

La variable de paquete `Parse::Eyapp::Node::PREFIXES` contiene la lista de prefijos de los nombres de tipo que serán eliminados cuando `str` muestra el término. Por defecto contiene la cadena `'Parse::Eyapp::Node::'`.

- `our $INDENT = 0` La variable `$Parse::Eyapp::Node::INDENT` controla el formato de presentación usado por `Parse::Eyapp::Node::str` :
  1. Si es 0 la representación es compacta.
  2. Si es 1 se usará un sangrado razonable.
  3. Si es 2 cada paréntesis cerrar que este a una distancia mayor de `$Parse::Eyapp::Node::LINESEP` líneas será comentado con el tipo del nodo. Por defecto la variable `$Parse::Eyapp::Node::LINESEP` está a 4. Véase la página 699 para un ejemplo con `$INDENT` a 2 y `$Parse::Eyapp::Node::LINESEP` a 4.

- `our $STRSEP = ','`

Separador de nodos en un término. Por defecto es la coma.

- `our $DELIMITER = '['`

Delimitador de presentación de la información proveída por el método `info`. Por defecto los delimitadores son corchetes. Puede elegirse uno cualquiera de la lista:

'[', '{', '(', '<'

si se define como la cadena vacía o `undef` no se usarán delimitadores.

- `our $FOOTNOTE_HEADER = "\n-----\n"`

Delimitador para las notas a pie de página.

- `our $FOOTNOTE_SEP = ")\n"` Separador de la etiqueta/número de la nota al pie

- `our $FOOTNOTE_LEFT = '^{'`, `our $FOOTNOTE_RIGHT = '}'`

Definen el texto que rodea al número de referencia en el nodo.

## Notas a Pié de Página o Footnotes

El siguiente ejemplo ilustra el manejo de *notas a pié de árbol* usando `str`. Se han definido los siguientes métodos:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> sed -ne '677,$p' Simple6.eyp
$Parse::Eyapp::Node::INDENT = 1;
sub TERMINAL::info {
 my @a = join ':', @{$_[0]->{attr}};
 return "@a"
}

sub PROGRAM::footnote {
 return "Types:\n"
 .Dumper($_[0]->{types}).
 "Symbol Table:\n"
 .Dumper($_[0]->{symboltable})
}

sub FUNCTION::info {
 return $_[0]->{function_name}[0]
}

sub FUNCTION::footnote {
 return Dumper($_[0]->{symboltable})
}

sub BLOCK::info {
 return $_[0]->{line}
}

sub VAR::info {
 return $_[0]->{definition}{type} if defined $_[0]->{definition}{type};
 return "No declarado!";
}

*FUNCTIONCALL::info = *VARARRAY::info = \&VAR::info;
```

El resultado para el programa de entrada:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> cat salida
int a,b;

int f(char c) {
```

```

a[2] = 4;
b[1][3] = a + b;
c = c[5] * 2;
return g(c);
}

```

produce una salida por `stderr`:

Identifier g not declared

y la siguiente descripción de la estructura de datos:

```

PROGRAM^{0}(
 FUNCTION{f}^{1}(
 ASSIGN(
 VARARRAY{INT}(
 TERMINAL{a:4},
 INDEXSPEC(
 INUM(
 TERMINAL{2:4}
)
)
),
 INUM(
 TERMINAL{4:4}
)
),
 ASSIGN(
 VARARRAY{INT}(
 TERMINAL{b:5},
 INDEXSPEC(
 INUM(
 TERMINAL{1:5}
),
 INUM(
 TERMINAL{3:5}
)
)
),
 PLUS(
 VAR{INT}(
 TERMINAL{a:5}
),
 VAR{INT}(
 TERMINAL{b:5}
)
)
),
 ASSIGN(
 VAR{CHAR}(
 TERMINAL{c:6}
),
 TIMES(
 VARARRAY{CHAR}(
 TERMINAL{c:6},

```



```

 INDEXSPEC(
 INUM(
 TERMINAL{5:6}
)
)
),
 INUM(
 TERMINAL{2:6}
)
),
RETURN(
 FUNCTIONCALL{No declarado!}(
 TERMINAL{g:7},
 ARGLIST(
 VAR{CHAR}(
 TERMINAL{c:7}
)
)
)
)

```

0)

\$VAR1 = {

} ;

\$VAR1 = {

```

},
'b' => {
 'type' => 'INT',
 'line' => 1
},
'f' => {
 'type' => 'F(X_1(CHAR),INT)',
 'line' => 3
}
};

```

```

1)
$VAR1 = {
 'c' => {
 'type' => 'CHAR',
 'param' => 1,
 'line' => 3
 }
};

```

### Usando str sobre una Lista de Árboles

Para usar **str** sobre una lista de árboles la llamada deberá hacerse como método de clase, i.e. con el nombre de la clase como prefijo y no con el objeto:

```
Parse::Eyapp::Node->str(@forest)
```

cuando quiera convertir a su representación término mas de un árbol. El código:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code> sed -ne '652,658p' Simple4.eypp \
| cat -n
1 local $Parse::Eyapp::Node::INDENT = 2;
2 local $Parse::Eyapp::Node::DELIMITER = "";
3 print $t->str."\n";
4 {
5 local $" = "\n";
6 print Parse::Eyapp::Node->str(@blocks)."\n";
7 }

```

produce la salida:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code> eyapp Simple4 ;\
treereg SimpleTrans.trg ;\
usesimple4.pl

```

\*\*\*\*\*

```

f() {
 int a,b[1][2],c[1][2][3];
 char d[10];
 b[0][1] = a;
}

```

```

PROGRAM(
 FUNCTION[f](
 ASSIGN(

```

```

VARARRAY(
 TERMINAL[b:4],
 INDEXSPEC(
 INUM(
 TERMINAL[0:4]
),
 INUM(
 TERMINAL[1:4]
)
) # INDEXSPEC
) # VARARRAY,
VAR(
 TERMINAL[a:4]
)
) # ASSIGN
) # FUNCTION
) # PROGRAM

Match[PROGRAM:0:blocks](
 Match[FUNCTION:1:blocks:[f]]
)

```

Obsérvense los comentarios `# TIPODENODO` que acompañan a los paréntesis cerrar cuando estos están lejos (esto es, a una distancia de mas de `$Parse::Eyapp::Node::LINESEP` líneas) del correspondiente paréntesis abrir. Tales comentarios son consecuencia de haber establecido el valor de `$Parse::Eyapp::Node::INDENT` a 2.

## 8.22. El Método descendant

El método `descendant` es un método de los objetos `Parse::Eyapp::Node` y retorna una referencia al nodo descrito por la cadena de coordenadas que se le pasa como argumento. En este sentido el método `child` es una especialización de `descendant`.

```

DB<7> x $t->child(0)->child(0)->child(1)->child(0)->child(2)->child(1)->str
0 ,
BLOCK[8:4:test]~{0}(
 CONTINUE[10,10]
)
DB<8> x $t->descendant('0.0.1.0.2.1')->str
0 ,
BLOCK[8:4:test]~{0}(
 CONTINUE[10,10]
)

```

## 8.23. Conceptos Básicos del Análisis LR

Los analizadores generados por `eyapp` entran en la categoría de analizadores *LR*. Estos analizadores construyen una derivación a derechas inversa (o *antiderivación*). De ahí la R en LR (del inglés *rightmost derivation*). El árbol sintáctico es construido de las hojas hacia la raíz, siendo el último paso en la antiderivación la construcción de la primera derivación desde el símbolo de arranque.

Empezaremos entonces considerando las frases que pueden aparecer en una derivación a derechas. Tales frases consituyen el lenguaje *FSD*:

**Definición 8.23.1.** Dada una gramática  $G = (\Sigma, V, P, S)$  no ambigua, se denota por  $FSD$  (lenguaje de las formas Sentenciales a Derechas) al lenguaje de las sentencias que aparecen en una derivación a derechas desde el símbolo de arranque.

$$FSD = \left\{ \alpha \in (\Sigma \cup V)^* : \exists S \xRightarrow[RM]{*} \alpha \right\}$$

Donde la notación  $RM$  indica una derivación a derechas (rightmost). Los elementos de  $FSD$  se llaman “formas sentenciales derechas”.

Dada una gramática no ambigua  $G = (\Sigma, V, P, S)$  y una frase  $x \in L(G)$  el proceso de antiderivación consiste en encontrar la última derivación a derechas que dió lugar a  $x$ . Esto es, si  $x \in L(G)$  es porque existe una derivación a derechas de la forma

$$S \xRightarrow{*} yAz \Rightarrow ywz = x.$$

El problema es averiguar que regla  $A \rightarrow w$  se aplicó y en que lugar de la cadena  $x$  se aplicó. En general, si queremos antiderivar una forma sentencial derecha  $\beta\alpha w$  debemos averiguar por que regla  $A \rightarrow \alpha$  seguir y en que lugar de la forma (después de  $\beta$  en el ejemplo) aplicarla.

$$S \xRightarrow{*} \beta Aw \Rightarrow \beta\alpha w.$$

La pareja formada por la regla y la posición se denomina mango o manecilla de la forma. Esta denominación viene de la visualización gráfica de la regla de producción como una mano que nos permite escalar hacia arriba en el árbol. Los “dedos” serían los símbolos en la parte derecha de la regla de producción.

**Definición 8.23.2.** Dada una gramática  $G = (\Sigma, V, P, S)$  no ambigua, y dada una forma sentencial derecha  $\alpha = \beta\gamma x$ , con  $x \in \Sigma^*$ , el mango o handle de  $\alpha$  es la última producción/posición que dió lugar a  $\alpha$ :

$$\xRightarrow[RM]{*} S \Rightarrow \beta Bx \Rightarrow \beta\gamma x = \alpha$$

Escribiremos:  $handle(\alpha) = (B \rightarrow \gamma, \beta\gamma)$ . La función  $handle$  tiene dos componentes:  $handle_1(\alpha) = B \rightarrow \gamma$  y  $handle_2(\alpha) = \beta\gamma$

Si dispusiéramos de un procedimiento que fuera capaz de identificar el mango, esto es, de detectar la regla y el lugar en el que se posiciona, tendríamos un mecanismo para construir un analizador. Lo curioso es que, a menudo es posible encontrar un autómata finito que reconoce el lenguaje de los prefijos  $\beta\gamma$  que terminan en el mango. Con mas precisión, del lenguaje:

**Definición 8.23.3.** El conjunto de prefijos viables de una gramática  $G$  se define como el conjunto:

$$PV = \left\{ \delta \in (\Sigma \cup V)^* : \exists S \xRightarrow[RM]{*} \alpha \text{ y } \delta \text{ es un prefijo de } handle_2(\alpha) \right\}$$

Esto es, el lenguaje de los prefijos viables es el conjunto de frases que son prefijos de  $handle_2(\alpha) = \beta\gamma$ , siendo  $\alpha$  una forma sentencial derecha ( $\alpha \in FSD$ ). Los elementos de  $PV$  se denominan prefijos viables.

Obsérvese que si se dispone de un autómata que reconoce  $PV$  entonces se dispone de un mecanismo para investigar el lugar y el aspecto que pueda tener el mango. Si damos como entrada la sentencia  $\alpha = \beta\gamma x$  a dicho autómata, el autómata aceptará la cadena  $\beta\gamma$  pero rechazará cualquier extensión del prefijo. Ahora sabemos que el mango será alguna regla de producción de  $G$  cuya parte derecha sea un sufijo de  $\beta\gamma$ .

**Definición 8.23.4.** El siguiente autómata finito no determinista puede ser utilizado para reconocer el lenguaje de los prefijos viables PV:

- Alfabeto =  $V \cup \Sigma$
- Los estados del autómata se denominan  $LR(0)$  items. Son parejas formadas por una regla de producción de la gramática y una posición en la parte derecha de la regla de producción. Por ejemplo,  $(E \rightarrow E + E, 2)$  sería un  $LR(0)$  item para la gramática de las expresiones.

Conjunto de Estados:

$$Q = \{(A \rightarrow \alpha, n) : A \rightarrow \alpha \in P, n \leq |\alpha|\}$$

La notación  $|\alpha|$  denota la longitud de la cadena  $\alpha$ . En vez de la notación  $(A \rightarrow \alpha, n)$  escribiremos:  $A \rightarrow \beta \uparrow \gamma = \alpha$ , donde la flecha ocupa el lugar indicado por el número  $n = |\beta|$ :

- Función de transición:

$$\delta(A \rightarrow \alpha \uparrow X \beta, X) = A \rightarrow \alpha X \uparrow \beta \quad \forall X \in V \cup \Sigma$$

$$\delta(A \rightarrow \alpha \uparrow B \beta, \epsilon) = B \rightarrow \gamma \uparrow B \quad \forall B \in V$$

- Estado de arranque: Se añade la “superregla”  $S' \rightarrow S$  a la gramática  $G = (\Sigma, V, P, S)$ . El  $LR(0)$  item  $S' \rightarrow \uparrow S$  es el estado de arranque.
- Todos los estados definidos (salvo el de muerte) son de aceptación.

Denotaremos por  $LR(0)$  a este autómata. Sus estados se denominan  $LR(0)$  – items. La idea es que este autómata nos ayuda a reconocer los prefijos viables PV.

Una vez que se tiene un autómata que reconoce los prefijos viables es posible construir un analizador sintáctico que construye una antiderivación a derechas. La estrategia consiste en “alimentar” el autómata con la forma sentencial derecha. El lugar en el que el autómata se detiene, rechazando indica el lugar exacto en el que termina el *handle* de dicha forma.

**Ejemplo 8.23.1.** Consideremos la gramática:

$$\begin{aligned} S &\rightarrow a S b \\ S &\rightarrow \epsilon \end{aligned}$$

El lenguaje generado por esta gramática es  $L(G) = \{a^n b^n : n \geq 0\}$ . Es bien sabido que el lenguaje  $L(G)$  no es regular. La figura 8.23.1 muestra el autómata finito no determinista con  $\epsilon$ -transiciones (NFA) que reconoce los prefijos viables de esta gramática, construido de acuerdo con el algoritmo 8.23.4.

**Ejercicio 8.23.1.** Simule el comportamiento del autómata sobre la entrada *aabb*. ¿Donde rechaza? ¿En que estados está el autómata en el momento del rechazo?. ¿Qué etiquetas tienen? Haga también las trazas del autómata para las entradas *aaSbb* y *aSb*. ¿Que antiderivación ha construido el autómata con sus sucesivos rechazos? ¿Que terminales se puede esperar que hayan en la entrada cuando se produce el rechazo del autómata?

## 8.24. Construcción de las Tablas para el Análisis SLR

### 8.24.1. Los conjuntos de Primeros y Siguietes

Repasemos las nociones de conjuntos de *Primeros* y *siguietes*:

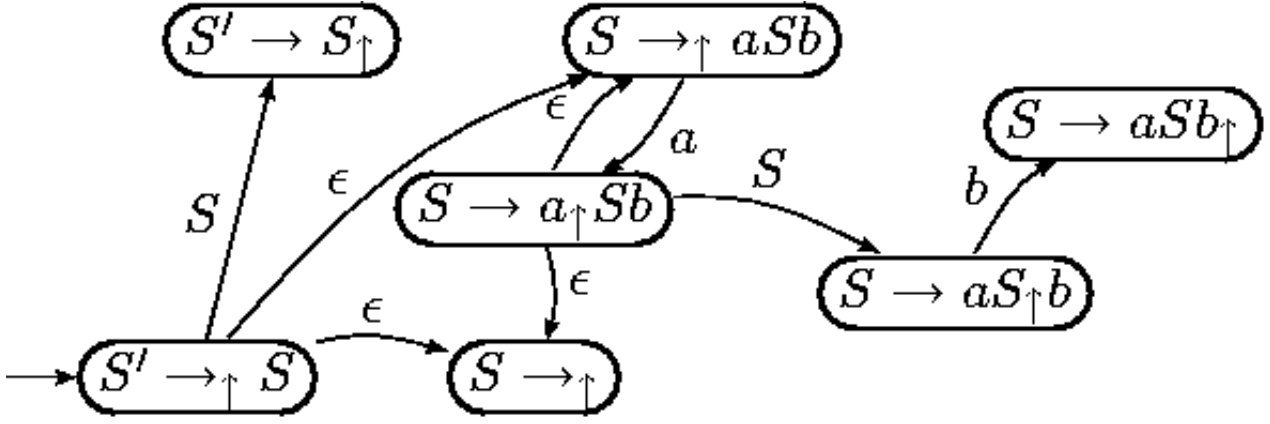


Figura 8.1: NFA que reconoce los prefijos viables

**Definición 8.24.1.** Dada una gramática  $G = (\Sigma, V, P, S)$  y un símbolo  $\alpha \in (V \cup \Sigma)^*$  se define el conjunto  $FIRST(\alpha)$  como:

$$FIRST(\alpha) = \{b \in \Sigma : \alpha \xRightarrow{*} b\beta\} \cup N(\alpha)$$

donde:

$$N(\alpha) = \begin{cases} \{\epsilon\} & \text{si } \alpha \xRightarrow{*} \epsilon \\ \emptyset & \text{en otro caso} \end{cases}$$

**Definición 8.24.2.** Dada una gramática  $G = (\Sigma, V, P, S)$  y una variable  $A \in V$  se define el conjunto  $FOLLOW(A)$  como:

$$FOLLOW(A) = \{b \in \Sigma : \exists S \xRightarrow{*} \alpha Ab\beta\} \cup E(A)$$

donde

$$E(A) = \begin{cases} \{\$ \} & \text{si } S \xRightarrow{*} \alpha A \\ \emptyset & \text{en otro caso} \end{cases}$$

**Algoritmo 8.24.1.** Construcción de los conjuntos  $FIRST(X)$

1. Si  $X \in \Sigma$  entonces  $FIRST(X) = X$
2. Si  $X \rightarrow \epsilon$  entonces  $FIRST(X) = FIRST(X) \cup \{\epsilon\}$
3. Si  $X \in V$  y  $X \rightarrow Y_1 Y_2 \dots Y_k \in P$  entonces

$i = 1;$

do

$FIRST(X) = FIRST(X) \cup FIRST(Y_i) - \{\epsilon\};$

$i++;$

mientras  $(\epsilon \in FIRST(Y_i) \text{ and } (i \leq k))$

si  $(\epsilon \in FIRST(Y_k) \text{ and } i > k)$   $FIRST(X) = FIRST(X) \cup \{\epsilon\}$

Este algoritmo puede ser extendido para calcular  $FIRST(\alpha)$  para  $\alpha = X_1 X_2 \dots X_n \in (V \cup \Sigma)^*$ .

**Algoritmo 8.24.2.** Construcción del conjunto  $FIRST(\alpha)$

$i = 1;$

$FIRST(\alpha) = \emptyset;$

do

$FIRST(\alpha) = FIRST(\alpha) \cup FIRST(X_i) - \{\epsilon\};$

$i++;$

mientras  $(\epsilon \in FIRST(X_i) \text{ and } (i \leq n))$

si  $(\epsilon \in FIRST(X_n) \text{ and } i > n)$   $FIRST(\alpha) = FIRST(X) \cup \{\epsilon\}$

**Algoritmo 8.24.3.** Construcción de los conjuntos  $FOLLOW(A)$  para las variables sintácticas  $A \in V$ :  
Repetir los siguientes pasos hasta que ninguno de los conjuntos  $FOLLOW$  cambie:

1.  $FOLLOW(S) = \{\$ \}$  ( $\$$  representa el final de la entrada)
2. Si  $A \rightarrow \alpha B \beta$  entonces

$$FOLLOW(B) = FOLLOW(B) \cup (FIRST(\beta) - \{\epsilon\})$$

3. Si  $A \rightarrow \alpha B$  o bien  $A \rightarrow \alpha B \beta$  y  $\epsilon \in FIRST(\beta)$  entonces

$$FOLLOW(B) = FOLLOW(B) \cup FOLLOW(A)$$

### 8.24.2. Construcción de las Tablas

Para la construcción de las tablas de un analizador SLR se construye el *autómata finito determinista* (DFA)  $(Q, \Sigma, \delta, q_0)$  equivalente al NFA presentado en la sección 8.23 usando el *algoritmo de construcción del subconjunto*.

Como recordará, en la construcción del subconjunto, partiendo del estado de arranque  $q_0$  del NFA con  $\epsilon$ -transiciones se calcula su *clausura*  $\overline{\{q_0\}}$  y las clausuras de los conjuntos de estados  $\delta(\overline{\{q_0\}}, a)$  a los que transita. Se repite el proceso con los conjuntos resultantes hasta que no se introducen nuevos conjuntos-estado.

#### Clausura de un Conjunto de Estados

La clausura  $\overline{A}$  de un subconjunto de estados del autómata  $A$  esta formada por todos los estados que pueden ser alcanzados mediante transiciones etiquetadas con la palabra vacía (denominadas  $\epsilon$  transiciones) desde los estados de  $A$ . Se incluyen en  $\overline{A}$ , naturalmente los estados de  $A$ .

$$\overline{A} = \{q \in Q \mid \exists q' \in A : \hat{\delta}(q, \epsilon) = q'\}$$

Aquí  $\hat{\delta}$  denota la *función de transición del autómata* extendida a cadenas de  $\Sigma^*$ .

$$\hat{\delta}(q, x) = \begin{cases} \delta(\hat{\delta}(q, y), a) & \text{si } x = ya \\ q & \text{si } x = \epsilon \end{cases} \quad (8.1)$$

En la práctica, y a partir de ahora así lo haremos, se prescinde de diferenciar entre  $\delta$  y  $\hat{\delta}$  usándose indistintamente la notación  $\delta$  para ambas funciones.

La clausura puede ser computada usando una estructura de pila o aplicando la expresión recursiva dada en la ecuación 8.1. Una forma de computarla viene dada por el siguiente pseudocódigo:

```
function closure(I : set of LR(0)-items)
begin
 J = I;
 repeat
 changes = FALSE;
 for A->alpha . B beta in J do
 for B->gamma in G do
 next if B->.gamma in J
 insert B->.gamma in J
 changes = TRUE;
 end for
 end for
 until nochanges;
 return J
end
```

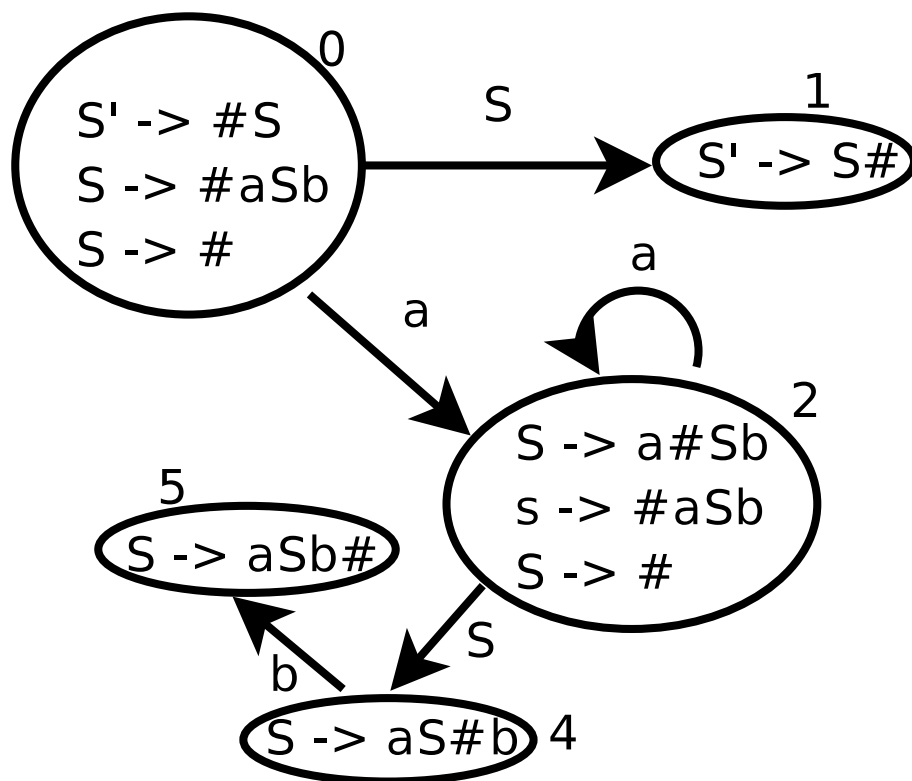


Figura 8.2: DFA equivalente al NFA de la figura 8.23.1

### Ejemplo de Construcción del DFA

Para el NFA mostrado en el ejemplo 8.23.1 el DFA construido mediante esta técnica es el que se muestra en la figura 8.24.2. Se ha utilizado el símbolo  $\#$  como marcador. Se ha omitido el número 3 para que los estados coincidan en numeración con los generados por **eyapp** (véase el cuadro 8.24.2).

### Las Tablas de Saltos y de Acciones

Un analizador sintáctico LR utiliza una tabla para su análisis. Esa tabla se construye a partir de la tabla de transiciones del DFA. De hecho, la tabla se divide en dos tablas, una llamada *tabla de saltos* o *tabla de gotos* y la otra *tabla de acciones*.

La *tabla goto* de un analizador *SLR* no es más que la tabla de transiciones del autómata DFA obtenido aplicando la construcción del subconjunto al NFA definido en 8.23.4. De hecho es la tabla de transiciones restringida a  $V$  (recuerde que el alfabeto del autómata es  $V \cup \Sigma$ ). Esto es,

$$\delta_{|V \times Q} : V \times Q \rightarrow Q.$$

donde se define  $goto(i, A) = \delta(A, I_i)$

La parte de la función de transiciones del DFA que corresponde a los terminales que no producen rechazo, esto es,  $\delta_{|\Sigma \times Q} : \Sigma \times Q \rightarrow Q$  se adjunta a una tabla que se denomina *tabla de acciones*. La tabla de acciones es una tabla de doble entrada en los estados y en los símbolos de  $\Sigma$ . Las acciones de transición ante terminales se denominan *acciones de desplazamiento* o (*acciones shift*):

$$\delta_{|\Sigma \times Q} : \Sigma \times Q \rightarrow Q$$

donde se define  $action(i, a) = \delta(a, I_i)$

### Ante que Terminales se debe Reducir

Cuando un estado  $s$  contiene un LR(0)-item de la forma  $A \rightarrow \alpha \uparrow$ , esto es, el estado corresponde a un posible rechazo, ello indica que hemos llegado a un final del prefijo viable, que hemos visto  $\alpha$  y



que, por tanto, es probable que  $A \rightarrow \alpha$  sea el *handle* de la forma sentencial derecha actual. Por tanto, añadiremos en entradas de la forma  $(s, a)$  de la tabla de acciones una acción que indique que hemos encontrado el mango en la posición actual y que la regla asociada es  $A \rightarrow \alpha$ . A una acción de este tipo se la denomina *acción de reducción*.

La cuestión es, ¿para que valores de  $a \in \Sigma$  debemos disponer que la acción para  $(s, a)$  es de reducción? Podríamos decidir que ante cualquier terminal  $a \in \Sigma$  que produzca un rechazo del autómata, pero podemos ser un poco mas selectivos. No cualquier terminal puede estar en la entrada en el momento en el que se produce la antiderivación o reducción. Observemos que si  $A \rightarrow \alpha$  es el *handle* de  $\gamma$  es porque:

$$\begin{array}{ccccc} & * & & * & \\ \exists S & \Longrightarrow & \beta A b x & \Longrightarrow & \beta \alpha b x = \gamma \\ & RM & & RM & \end{array}$$

Por tanto, cuando estamos reduciendo por  $A \rightarrow \alpha$  los únicos terminales legales que cabe esperar en una reducción por  $A \rightarrow \alpha$  son los terminales  $b \in FOLLOW(A)$ .

### Algoritmo de Construcción de Las Tablas SLR

Dada una gramática  $G = (\Sigma, V, P, S)$ , podemos construir las tablas de acciones (*action table*) y transiciones (*gotos table*) mediante el siguiente algoritmo:

#### Algoritmo 8.24.4. Construcción de Tablas SLR

1. Utilizando el Algoritmo de Construcción del Subconjunto, se construye el Autómata Finito Determinista (DFA)  $(Q, V \cup \Sigma, \delta, I_0, F)$  equivalente al Autómata Finito No Determinista (NFA) definido en 8.23.4. Sea  $C = \{I_1, I_2, \dots, I_n\}$  el conjunto de estados del DFA. Cada estado  $I_i$  es un conjunto de LR(0)-items o estados del NFA. Asociemos un índice  $i$  con cada conjunto  $I_i$ .
2. La tabla de gotos no es más que la función de transición del autómata restringida a las variables de la gramática:

$$goto(i, A) = \delta(I_i, A) \text{ para todo } A \in V$$

3. Las acciones para el estado  $I_i$  se determinan como sigue:

- a) Si  $A \rightarrow \alpha \uparrow a \beta \in I_i$ ,  $\delta(I_i, a) = I_j$ ,  $a \in \Sigma$  entonces:

$$action[i][a] = shift\ j$$

- b) Si  $S' \rightarrow S \uparrow \in I_i$  entonces

$$action[i][\$] = accept$$

- c) Para cualquier otro caso de la forma  $A \rightarrow \alpha \uparrow \in I_i$  distinto del anterior hacer

$$\forall a \in FOLLOW(A) : action[i][a] = reduce\ A \rightarrow \alpha$$

4. Las entradas de la tabla de acción que queden indefinidas después de aplicado el proceso anterior corresponden a acciones de “error”.

### Conflictos en Un Analizador SLR

**Definición 8.24.3.** Si alguna de las entradas de la tabla resulta multievaluada, decimos que existe un conflicto y que la gramática no es SLR.

1. En tal caso, si una de las acciones es de ‘reducción’ y la otra es de ‘desplazamiento’, decimos que hay un conflicto shift-reduce o conflicto de desplazamiento-reducción.
2. Si las dos reglas indican una acción de reducción, decimos que tenemos un conflicto reduce-reduce o de reducción-reducción.

## Ejemplo de Cálculo de las Tablas SLR

**Ejemplo 8.24.1.** Al aplicar el algoritmo 8.24.4 a la gramática 8.23.1

|   |                          |
|---|--------------------------|
| 1 | $S \rightarrow a S b$    |
| 2 | $S \rightarrow \epsilon$ |

partiendo del autómata finito determinista que se construyó en la figura 8.24.2 y calculando los conjuntos de primeros y siguientes

|          | <i>FIRST</i>          | <i>FOLLOW</i>   |
|----------|-----------------------|-----------------|
| <i>S</i> | <i>a</i> , $\epsilon$ | <i>b</i> , $\$$ |

obtenemos la siguiente tabla de acciones SLR:

|   | <i>a</i>  | <i>b</i>  | $\$$           |
|---|-----------|-----------|----------------|
| 0 | <i>s2</i> | <i>r2</i> | <i>r2</i>      |
| 1 |           |           | <i>aceptar</i> |
| 2 | <i>s2</i> | <i>r2</i> | <i>r2</i>      |
| 4 |           | <i>s5</i> |                |
| 5 |           | <i>r1</i> | <i>r1</i>      |

Las entradas denotadas con *s n* (*s* por *shift*) indican un desplazamiento al estado *n*, las denotadas con *r n* (*r* por *reduce* o *reducción*) indican una operación de reducción o antiderivación por la regla *n*. Las entradas vacías corresponden a acciones de error.

### Las Tablas Construidas por **eyapp**

El método de análisis *LALR* usado por **eyapp** es una extensión del método SLR esbozado aquí. Supone un compromiso entre potencia (conjunto de gramáticas englobadas) y eficiencia (cantidad de memoria utilizada, tiempo de proceso). Veamos como **eyapp** aplica la construcción del subconjunto a la gramática del ejemplo 8.23.1. Para ello construimos el siguiente programa **eyapp**:

```
$ cat -n aSb.y
1 %%
2 S: # empty
3 | 'a' S 'b'
4 ;
5 %%
.....
```

y compilamos, haciendo uso de la opción **-v** para que **eyapp** produzca las tablas en el fichero **aSb.output**:

```
$ ls -l aSb.*
-rw-r--r-- 1 lhp lhp 738 2004-12-19 09:52 aSb.output
-rw-r--r-- 1 lhp lhp 1841 2004-12-19 09:52 aSb.pm
-rw-r--r-- 1 lhp lhp 677 2004-12-19 09:46 aSb.y
```

El contenido del fichero **aSb.output** se muestra en la tabla 8.24.2. Los números de referencia a las producciones en las acciones de reducción vienen dados por:

```
0: $start -> S $end
1: S -> /* empty */
2: S -> 'a' S 'b'
```

Observe que el final de la entrada se denota por **\$end** y el marcador en un LR-item por un punto. Fíjese en el estado 2: En ese estado están también los items

$S \rightarrow . 'a' S 'b' \text{ y } S \rightarrow .$

sin embargo no se explicitan por que se entiende que su pertenencia es consecuencia directa de aplicar la operación de clausura. Los LR items cuyo marcador no está al principio se denominan *items núcleo*.

| Estado 0                                                                                      | Estado 1                                         | Estado 2                                                                                  |
|-----------------------------------------------------------------------------------------------|--------------------------------------------------|-------------------------------------------------------------------------------------------|
| $\$start \rightarrow . S \$end$<br>'a'shift 2<br>\$default<br>reduce 1 (S)<br>S go to state 1 | $\$start \rightarrow S . \$end$<br>\$end shift 3 | $S \rightarrow 'a' . S 'b'$<br>'a'shift 2<br>\$default<br>reduce 1 (S)<br>S go to state 4 |
| Estado 3                                                                                      | Estado 4                                         | Estado 5                                                                                  |
| $\$start \rightarrow S \$end .$<br>\$default accept                                           | $S \rightarrow 'a' S . 'b'$<br>'b'shift 5        | $S \rightarrow 'a' S 'b' .$<br>\$default<br>reduce 2 (S)                                  |

Cuadro 8.2: Tablas generadas por `eyapp`. El estado 3 resulta de transitar con \$

Puede encontrar el listado completo de las tablas en `aSb.output` en el apéndice que se encuentra en la página ??.

**Ejercicio 8.24.1.** Compare la tabla 8.24.2 resultante de aplicar `eyapp` con la que obtuvo en el ejemplo 8.24.1.

## 8.25. Algoritmo de Análisis LR

Asi pues la tabla de transiciones del autómata nos genera dos tablas: la tabla de acciones y la de saltos. El algoritmo de análisis sintáctico *LR* en el que se basa `eyapp` utiliza una pila y dos tablas para analizar la entrada. Como se ha visto, la tabla de acciones contiene cuatro tipo de acciones:

1. Desplazar (*shift*)
2. Reducir (*reduce*)
3. Aceptar
4. Error

El algoritmo utiliza una pila en la que se guardan los estados del autómata. De este modo se evita tener que “comenzar” el procesado de la forma sentencial derecha resultante después de una reducción (antiderivación).

### Algoritmo 8.25.1. Analizador LR

```

my $parse = shift;
my @stack;
my $s0 = $parse->startstate;
push(@stack, $s0);
my $b = $parse->yylex();

```

```

FOREVER: {
 my $s = top(0);
 my $a = $b;
 switch ($parse->action[$s->state][$a]) {
 case "shift t" :
 my $t;
 $t->{state} = t;
 $t->{attr} = $a->{attr};
 push($t);
 $b = $parse->yylex();
 break;
 case "reduce A ->alpha" :
 my $r;
 $r->{attr} = $parse->Semantic{A ->alpha}->(top(|alpha|-1)->attr, ... , top(0)->attr);
 pop(|alpha|); # Saquemos length(alpha) elementos de la pila
 $r->{state} = $parse->goto[top(0)][A];
 push($r);
 break;
 case "accept" : return ($s->attr);
 default : $parse->yyerror("syntax error");
 }
 redo FOREVER;
}

```

Como es habitual,  $|x|$  denota la longitud de la cadena  $x$ . La función `top(k)` devuelve el elemento que ocupa la posición  $k$  desde el *top* de la pila (esto es, está a profundidad  $k$ ). La función `pop(k)` extrae  $k$  elementos de la pila. La notación `state->attr` hace referencia al atributo asociado con cada estado. Denotamos por `$Semantic{A->alpha}` el código de la acción asociada con la regla  $A \rightarrow \alpha$ .

Todos los analizadores LR comparten, salvo pequeñas excepciones, el mismo algoritmo de análisis. Lo que más los diferencia es la forma en la que construyen las tablas. En **eyapp** la construcción de las tablas de *acciones* y *gotos* se realiza mediante el algoritmo *LALR*.

|    |                                                       |                                                        |
|----|-------------------------------------------------------|--------------------------------------------------------|
| 1  | pl@nereida:~/LEyapp/examples\$                        | pl@nereida:~/LEyapp/examples\$ cat -n aSb.output       |
| 2  | -----                                                 | 1--Rules:                                              |
| 3  | In state 0:                                           | 2 -----                                                |
| 4  | Stack:[0]                                             | 3 0:        \$start -> S \$end                         |
| 5  | aabb                                                  | 4 1:        S -> /* empty */                           |
| 6  | Need token. Got >a<                                   | 5 2:        S -> 'a' S 'b'                             |
| 7  | Shift and go to state 2.                              | 6 -----                                                |
| 8  | -----                                                 | 7--States:                                             |
| 9  | In state 2:                                           | 8 -----                                                |
| 10 | Stack:[0,2]                                           | 9 State 0:                                             |
| 11 | Need token. Got >a<                                   | 10                                                     |
| 12 | Shift and go to state 2.                              | 11            \$start -> . S \$end        (Rule 0)     |
| 13 | -----                                                 | 12-----                                                |
| 14 | In state 2:                                           | 13            'a'        shift, and go to state 2      |
| 15 | Stack:[0,2,2]                                         | 14                                                     |
| 16 | Need token. Got >b<                                   | 15            \$default        reduce using rule 1 (S) |
| 17 | Reduce using rule 1 (S --> /* empty */): S -> epsilon | 16            S        go to state 1                   |
| 18 | Back to state 2, then go to state 4.                  | 17-----                                                |
| 19 | -----                                                 | 18-----                                                |
| 20 | In state 4:                                           | 19 State 1:                                            |
| 21 | Stack:[0,2,2,4]                                       | 20                                                     |
| 22 | Shift and go to state 5.                              | 21            \$start -> S . \$end        (Rule 0)     |
| 23 | -----                                                 | 22-----                                                |
| 24 | In state 5:                                           | 23            \$end        shift, and go to state 3    |
| 25 | Stack:[0,2,2,4,5]                                     | 24                                                     |
| 26 | Don't need token.                                     | 25 State 2:                                            |
| 27 | Reduce using rule 2 (S --> a S b)                     | 26            S -> a S b                               |
| 28 | Back to state 2, then go to state 4.                  | 27            S -> 'a' . S 'b'        (Rule 2)         |
| 29 | -----                                                 | 28-----                                                |
| 30 | In state 4:                                           | 29            'a'        shift, and go to state 2      |
| 31 | Stack:[0,2,4]                                         | 30                                                     |
| 32 | Need token. Got >b<                                   | 31            \$default        reduce using rule 1 (S) |
| 33 | Shift and go to state 5.                              | 32                                                     |
| 34 | -----                                                 | 33----- S        go to state 4                         |
| 35 | In state 5:                                           | 34                                                     |
| 36 | Stack:[0,2,4,5]                                       | 35 State 3:                                            |
| 37 | Don't need token.                                     | 36                                                     |
| 38 | Reduce using rule 2 (S --> a S b)                     | 37            S -> a S b                               |
| 39 | Back to state 0, then go to state 1.                  | 38            \$start -> S \$end .        (Rule 0)     |
| 40 | -----                                                 | 39----- \$default        accept                        |
| 41 | In state 1:                                           | 40                                                     |
| 42 | Stack:[0,1]                                           | 41 State 4:                                            |
| 43 | Need token. Got ><                                    | 42                                                     |
| 44 | Shift and go to state 3.                              | 43            S -> 'a' S . 'b'        (Rule 2)         |
| 45 | -----                                                 | 44-----                                                |
| 46 | In state 3:                                           | 45            'b'        shift, and go to state 5      |
| 47 | Stack:[0,1,3]                                         | 46                                                     |
| 48 | Don't need token.                                     | 47 State 5:                                            |
| 49 | Accept.                                               | 48                                                     |
|    |                                                       | 49            S -> 'a' S 'b' .        (Rule 2)         |
|    |                                                       | 50                                                     |
|    |                                                       | 51            \$default        reduce using rule 2 (S) |
|    |                                                       | 52                                                     |
|    |                                                       | 53                                                     |
|    |                                                       | 54 Summary:                                            |
|    |                                                       | 55 -----                                               |
|    |                                                       | 56 Number of rules        : 3                          |

## 8.26. Precedencia y Asociatividad

Recordemos que si al construir la tabla LALR, alguna de las entradas de la tabla resulta multievaluada, decimos que existe un conflicto. Si una de las acciones es de ‘reducción’ y la otra es de ‘desplazamiento’, se dice que hay un *conflicto shift-reduce* o *conflicto de desplazamiento-reducción*. Si las dos reglas indican una acción de reducción, decimos que tenemos un *conflicto reduce-reduce* o de *reducción-reducción*. En caso de que no existan indicaciones específicas *eyapp* resuelve los conflictos que aparecen en la construcción de la tabla utilizando las siguientes reglas:

1. Un conflicto *reduce-reduce* se resuelve eligiendo la producción que se listó primero en la especificación de la gramática.
2. Un conflicto *shift-reduce* se resuelve siempre en favor del *shift*

Las declaraciones de precedencia y asociatividad mediante las palabras reservadas `%left` , `%right` , `%nonassoc` se utilizan para modificar estos criterios por defecto. La declaración de `token s` mediante la palabra reservada `%token` no modifica la precedencia. Si lo hacen las declaraciones realizadas usando las palabras `left` , `right` y `nonassoc` .

1. Los *tokens* declarados en la misma línea tienen igual precedencia e igual asociatividad. La precedencia es mayor cuanto mas abajo su posición en el texto. Así, en el ejemplo de la calculadora en la sección 7.1, el *token* `*` tiene mayor precedencia que `+` pero la misma que `/`.
2. La precedencia de una regla  $A \rightarrow \alpha$  se define como la del terminal mas a la derecha que aparece en  $\alpha$ . En el ejemplo, la producción

`expr : expr '+' expr`

tiene la precedencia del *token* `+`.

3. Para decidir en un conflicto *shift-reduce* se comparan la precedencia de la regla con la del terminal que va a ser desplazado. Si la de la regla es mayor se reduce si la del *token* es mayor, se desplaza.
4. Si en un conflicto *shift-reduce* ambos la regla y el terminal que va a ser desplazado tiene la misma precedencia *eyapp* considera la asociatividad, si es asociativa a izquierdas, reduce y si es asociativa a derechas desplaza. Si no es asociativa, genera un mensaje de error.  
Obsérvese que, en esta situación, la asociatividad de la regla y la del *token* han de ser por fuerza, las mismas. Ello es así, porque en *eyapp* los *tokens* con la misma precedencia se declaran en la misma línea y sólo se permite una declaración por línea.
5. *Por tanto es imposible declarar dos tokens con diferente asociatividad y la misma precedencia.*
6. Es posible modificar la precedencia “natural” de una regla, calificándola con un *token* específico. para ello se escribe a la derecha de la regla `prec token`, donde `token` es un *token* con la precedencia que deseamos. Vea el uso del *token dummy* en el siguiente ejercicio.

Para ilustrar las reglas anteriores usaremos el siguiente programa *eyapp*:

```
$ cat -n Precedencia.yp
 1 %token NUMBER
 2 %left '@'
 3 %right '&' dummy
 4 %%
 5 list
 6 :
 7 | list '\n'
```

```

8 | list e
9 ;
10
11 e : NUMBER
12 | e '&' e
13 | e '@' e %prec dummy
14 ;
15
16 %%

```

El código del programa cliente es el siguiente:

```

$ cat -n useprecedencia.pl
cat -n useprecedencia.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Precedencia;
4
5 sub Error {
6 exists $_[0]->YYData->{ERRMSG}
7 and do {
8 print $_[0]->YYData->{ERRMSG};
9 delete $_[0]->YYData->{ERRMSG};
10 return;
11 };
12 print "Syntax error.\n";
13 }
14
15 sub Lexer {
16 my($parser)=shift;
17
18 defined($parser->YYData->{INPUT})
19 or $parser->YYData->{INPUT} = <STDIN>
20 or return('','undef');
21
22 $parser->YYData->{INPUT} =~ s/^[\t]//;
23
24 for ($parser->YYData->{INPUT}) {
25 s/^[0-9]+(?:\.[0-9]+)?//
26 and return('NUMBER',$1);
27 s/^(.)//s
28 and return($1,$1);
29 }
30 }
31
32 my $debug_level = (@ARGV)? oct(shift @ARGV): 0x1F;
33 my $parser = Precedencia->new();
34 $parser->YYParse(yylex => \&Lexer, yyerror => \&Error, yydebug => $debug_level);

```

Observe la llamada al analizador en la línea 34. Hemos añadido el parámetro con nombre *yydebug* con argumento `yydebug => $debug_level` (véase la sección 8.3 para ver los posibles valores de depuración).

Compilamos a continuación el módulo usando la opción `-v` para producir información sobre los conflictos y las tablas de salto y de acciones:

```

eyapp -v -m Precedencia Precedencia.y
$ ls -ltr |tail -2
-rw-r--r-- 1 lhp lhp 1628 2004-12-07 13:21 Precedencia.pm
-rw-r--r-- 1 lhp lhp 1785 2004-12-07 13:21 Precedencia.output

```

La opción `-v` genera el fichero `Precedencia.output` el cual contiene información detallada sobre el autómata:

```

$ cat -n Precedencia.output
 1 Conflicts:
 2 -----
 3 Conflict in state 8 between rule 6 and token '@' resolved as reduce.
 4 Conflict in state 8 between rule 6 and token '&' resolved as shift.
 5 Conflict in state 9 between rule 5 and token '@' resolved as reduce.
 6 Conflict in state 9 between rule 5 and token '&' resolved as shift.
 7
 8 Rules:
 9 -----
10 0: $start -> list $end
11 1: list -> /* empty */
12 2: list -> list '\n'
13 3: list -> list e
14 4: e -> NUMBER
15 5: e -> e '&' e
16 6: e -> e '@' e
17 ...

```

¿Porqué se produce un conflicto en el estado 8 entre la regla 6 (`e -> e '@' e`) y el terminal '@'? Editando el fichero `Precedencia.output` podemos ver los contenidos del estado 8:

```

85 State 8:
86
87 e -> e . '&' e (Rule 5)
88 e -> e . '@' e (Rule 6)
89 e -> e '@' e . (Rule 6)
90
91 '&' shift, and go to state 7
92
93 $default reduce using rule 6 (e)

```

El ítem de la línea 88 indica que debemos desplazar, el de la línea 89 que debemos reducir por la regla 6. ¿Porqué `eyapp` resuelve el conflicto optando por reducir? ¿Que prioridad tiene la regla 6? ¿Que asociatividad tiene la regla 6? La declaración en la línea 13 modifica la precedencia y asociatividad de la regla:

```

13 | e '@' e %prec dummy

```

de manera que la regla pasa a tener la precedencia y asociatividad de `dummy`. Recuerde que habíamos declarado `dummy` como asociativo a derechas:

```

2 %left '@'
3 %right '&' dummy

```

¿Que ocurre? Que `dummy` tiene mayor prioridad que '@' y por tanto la regla tiene mayor prioridad que el terminal: por tanto se reduce.

¿Que ocurre cuando el terminal en conflicto es '&'? En ese caso la regla y el terminal tienen la misma prioridad. Se hace uso de la asociatividad a derechas que indica que el conflicto debe resolverse desplazando.



**Ejercicio 8.26.1.** *Explique la forma en que eyapp resuelve los conflictos que aparecen en el estado 9. Esta es la información sobre el estado 9:*

State 9:

```
e -> e . '&' e (Rule 5)
e -> e '&' e . (Rule 5)
e -> e . '@' e (Rule 6)
'&'shift, and go to state 7
$default reduce using rule 5 (e)
```

Veamos un ejemplo de ejecución:

```
$./useprecedencia.pl
```

```

In state 0:
Stack:[0]
Don't need token.
Reduce using rule 1 (list,0): Back to state 0, then go to state 1.
```

Lo primero que ocurre es una reducción por la regla en la que `list` produce vacío. Si miramos el estado 0 del autómata vemos que contiene:

```
20 State 0:
21
22 $start -> . list $end (Rule 0)
23
24 $default reduce using rule 1 (list)
25
26 list go to state 1
```

A continuación se transita desde 0 con `list` y se consume el primer terminal:

```

In state 1:
Stack:[0,1]
2@3@4
Need token. Got >NUMBER<
Shift and go to state 5.
```

```

In state 5:
Stack:[0,1,5]
Don't need token.
Reduce using rule 4 (e,1): Back to state 1, then go to state 2.
```

En el estado 5 se reduce por la regla `e -> NUMBER`. Esto hace que se retire el estado 5 de la pila y se transite desde el estado 1 viendo el símbolo `e`:

```

In state 2:
Stack:[0,1,2]
Need token. Got >@<
Shift and go to state 6.
```

```

In state 6:
Stack:[0,1,2,6]
```

Need token. Got >NUMBER<

Shift and go to state 5.

-----  
In state 5:

Stack: [0,1,2,6,5]

Don't need token.

Reduce using rule 4 (e,1): Back to state 6, then go to state 8.

-----  
In state 8:

Stack: [0,1,2,6,8]

Need token. Got >@<

Reduce using rule 6 (e,3): Back to state 1, then go to state 2.

-----  
...

Accept.

Obsérvese la resolución del conflicto en el estado 8

La presencia de conflictos, aunque no siempre, en muchos casos es debida a la introducción de ambigüedad en la gramática. Si el conflicto es de desplazamiento-reducción se puede resolver explicitando alguna regla que rompa la ambigüedad. Los conflictos de reducción-reducción suelen producirse por un diseño erróneo de la gramática. En tales casos, suele ser mas adecuado modificar la gramática.

## 8.27. Acciones en Medio de una Regla

A veces necesitamos insertar una acción en medio de una regla. Una acción en medio de una regla puede hacer referencia a los atributos de los símbolos que la preceden (usando \$n), pero no a los que le siguen.

Cuando se inserta una acción {action<sub>1</sub>} para su ejecución en medio de una regla  $A \rightarrow \alpha\beta$  :

$$A \rightarrow \alpha \{action_1\} \beta \{action_2\}$$

eyapp crea una variable sintáctica temporal  $T$  e introduce una nueva regla:

1.  $A \rightarrow \alpha T \beta \{action_2\}$

2.  $T \rightarrow \epsilon \{action_1\}$

Las acciones en mitad de una regla cuentan como un símbolo mas en la parte derecha de la regla. Asi pues, en una acción posterior en la regla, se deberán referenciar los atributos de los símbolos, teniendo en cuenta este hecho.

Las acciones en mitad de la regla pueden tener un atributo. Las acciones posteriores en la regla se referirán a él como \$<sub>n</sub>[n], siendo n su número de orden en la parte derecha.

Observe que la existencia de acciones intermedias implica que la gramática inicial es modificada. La introducción de las nuevas reglas puede dar lugar a ambigüedades y/o conflictos. Es responsabilidad del programador eliminarlos. Por ejemplo, dada la gramática:

```
cs : '{' decs ss '}' | '{' ss '}' ;
```

Si la modificamos como sigue:

```
cs : { decl(); } '{' decs ss '}'
 | '{' ss '}'
 ;
```

habremos introducido un conflicto.

**Ejercicio 8.27.1.** *Explique las razones por las cuales la introducción de la acción que prefija la primera regla da lugar a conflictos*

**Ejercicio 8.27.2.** *El conflicto no se arregla haciendo que la acción que precede sea la misma:*

```
cs : { decl(); } '{' decs ss '}'
 | { decl(); } '{' ss '}'
 ;
```

*Explique donde está el fallo de esta propuesta*

**Ejercicio 8.27.3.** *¿Se arregla el conflicto anterior usando esta otra alternativa?*

```
cs : tp '{' decs ss '}'
 | tp { decl(); } '{' ss '}'
 ;

tp : /* empty */ { decl(); }
 ;
```

## 8.28. Manejo en eyapp de Atributos Heredados

Supongamos que **eyapp** esta inmerso en la construcción de la antiderivación a derechas y que la forma sentencial derecha en ese momento es:

$$X_m \dots X_1 X_0 Y_1 \dots Y_n a_1 \dots a_0$$

y que el mango es  $B \rightarrow Y_1 \dots Y_n$  y en la entrada quedan por procesar  $a_1 \dots a_0$ .

Es posible acceder en **eyapp** a los valores de los atributos de los estados en la pila del analizador que se encuentran “por debajo” o si se quiere “a la izquierda” de los estados asociados con la regla por la que se reduce. Para ello se usa una llamada al método **YYSemval**. La llamada es de la forma  $\$_{[0]} \rightarrow \text{YYSemval}(\text{index})$ , donde **index** es un entero. Cuando se usan los valores  $1 \dots n$  devuelve lo mismo que  $\$_{[1]}, \dots \$_{[n]}$ . Esto es  $\$_{[1]}$  es el atributo asociado con  $Y_1$  y  $\$_{[n]}$  es el atributo asociado con  $Y_n$ . Cuando se usa con el valor 0 devolverá el valor del atributo asociado con el símbolo que esta a la izquierda del mango actual, esto es el atributo asociado con  $X_0$ , si se llama con -1 el que está dos unidades a la izquierda de la variable actual, esto es, el asociado con  $X_1$  etc. Así  $\$_{[-m]}$  denota el atributo de  $X_m$ .

Esta forma de acceder a los atributos es especialmente útil cuando se trabaja con *atributos heredados*. Esto es, cuando un atributo de un nodo del árbol sintáctico se computa en términos de valores de atributos de su padre y/o sus hermanos. Ejemplos de atributos heredados son la clase y tipo en la declaración de variables. Supongamos que tenemos el siguiente *esquema de traducción* para calcular la clase (C) y tipo (T) en las declaraciones (D) de listas (L) de identificadores:

|   |                 |                                                                                         |
|---|-----------------|-----------------------------------------------------------------------------------------|
| 1 | D $\rightarrow$ | C T { $\$L\{c\} = \$C\{c\}; \$L\{t\} = \$T\{t\}$ } L                                    |
| 2 | C $\rightarrow$ | global { $\$C\{c\} = \text{"global"}$ }                                                 |
| 3 | C $\rightarrow$ | local { $\$C\{c\} = \text{"local"}$ }                                                   |
| 4 | T $\rightarrow$ | integer { $\$T\{t\} = \text{"integer"}$ }                                               |
| 5 | T $\rightarrow$ | float { $\$T\{t\} = \text{"float"}$ }                                                   |
| 6 | L $\rightarrow$ | { $\$L_1\{t\} = \$L\{t\}; \$L_1\{c\} = \$L\{c\};$ } $L_1$ ','                           |
| 7 |                 | id { $\text{set\_class}(\$id\{v\}, \$L\{c\}); \text{set\_type}(\$id\{v\}, \$L\{t\});$ } |
| 8 | L $\rightarrow$ | id { $\text{set\_class}(\$id\{v\}, \$L\{c\}); \text{set\_type}(\$id\{v\}, \$L\{t\});$ } |

Los atributos **c** y **t** denotan respectivamente la clase y el tipo.

**Ejercicio 8.28.1.** *Evalúe el esquema de traducción para la entrada **global float x,y**. Represente el árbol de análisis, las acciones incrustadas y determine el orden de ejecución.*

*Olvide por un momento la notación usada en las acciones y suponga que se tratara de acciones **eyapp**. ¿En que orden construye **eyapp** el árbol y en que orden ejecutará las acciones?*

Observe que la simple razón por la cual la acción intermedia  $\{ \$L_1\{t\} = \$L\{t\}; \$L_1\{c\} = \$L\{c\}; \}$  no puede ser ejecutada en un programa **eyapp** es porque el nodo del árbol sintáctico asociado con  $L_1$  no existe en el momento en el que la acción es ejecutada. No es posible guardar algo en  $L_1\{c\}$  ya que no hay nodo asociado con  $L_1$ . La situación es distinta si se está trabajando con un esquema de traducción ya que estos primero construyen el árbol y luego lo recorren.

A la hora de transformar este esquema de traducción en un programa **eyapp** es importante darse cuenta que en cualquier derivación a derechas desde  $D$ , cuando se reduce por una de las reglas

$$L \rightarrow id \mid L_1 \text{ ',' } id$$

el símbolo a la izquierda de  $L$  es  $T$  y el que está a la izquierda de  $T$  es  $C$ . Considere, por ejemplo la derivación a derechas:

$$\begin{aligned} D &\Rightarrow C \ T \ L \Rightarrow C \ T \ L, id \Rightarrow C \ T \ L, id, id \Rightarrow C \ T \ id, id, id \Rightarrow \\ &\Rightarrow C \ float \ id, id, id \Rightarrow local \ float \ id, id, id \end{aligned}$$

Observe que el orden de recorrido de **eyapp** es:

$$\begin{aligned} local \ float \ id, id, id &\Leftarrow C \ float \ id, id \Leftarrow C \ T \ id, id, id \Leftarrow \\ &\Leftarrow C \ T \ L, id, id \Leftarrow C \ T \ L, id \Leftarrow C \ T \ L \Leftarrow D \end{aligned}$$

en la antiderivación, cuando el mango es una de las dos reglas para listas de identificadores,  $L \rightarrow id$  y  $L \rightarrow L, id$  es decir durante las tres últimas antiderivaciones:

$$C \ T \ L, id, id \Leftarrow C \ T \ L, id \Leftarrow C \ T \ L \Leftarrow D$$

las variables a la izquierda del mango son  $T$  y  $C$ . Esto ocurre siempre. Estas observaciones nos conducen al siguiente programa **eyapp**:

```
$ cat -n Inherited.y
1 %token FLOAT INTEGER
2 %token GLOBAL
3 %token LOCAL
4 %token NAME
5
6 %%
7 declarationlist
8 : # vacio
9 | declaration ';' declarationlist
10 ;
11
12 declaration
13 : class type namelist { ; }
14 ;
15
16 class
17 : GLOBAL
18 | LOCAL
19 ;
20
21 type
22 : FLOAT
23 | INTEGER
24 ;
25
```

```

26 namelist
27 : NAME
28 { printf("%s de clase %s, tipo %s\n",
29 $_[1], $_[0]->YYSemval(-1), $_[0]->YYSemval(0)); }
30 | namelist ',' NAME
31 { printf("%s de clase %s, tipo %s\n",
32 $_[3], $_[0]->YYSemval(-1), $_[0]->YYSemval(0)); }
33 ;
34 %%

```

A continuación escribimos el programa que usa el módulo generado por `eyapp`:

```

$ cat -n useinherited.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Inherited;
4
5 sub Error {
6 exists $_[0]->YYData->{ERRMSG}
7 and do {
8 print $_[0]->YYData->{ERRMSG};
9 delete $_[0]->YYData->{ERRMSG};
10 return;
11 };
12 print "Error sintáctico\n";
13 }
14
15 { # hagamos una clausura con la entrada
16 my $input;
17 local $/ = undef;
18 print "Entrada (En Unix, presione CTRL-D para terminar):\n";
19 $input = <stdin>;
20
21 sub scanner {
22
23 { # Con el redo del final hacemos un bucle "infinito"
24 if ($input =~ m|\G\s*INTEGER\b|igc) {
25 return ('INTEGER', 'INTEGER');
26 }
27 elsif ($input =~ m|\G\s*FLOAT\b|igc) {
28 return ('FLOAT', 'FLOAT');
29 }
30 elsif ($input =~ m|\G\s*LOCAL\b|igc) {
31 return ('LOCAL', 'LOCAL');
32 }
33 elsif ($input =~ m|\G\s*GLOBAL\b|igc) {
34 return ('GLOBAL', 'GLOBAL');
35 }
36 elsif ($input =~ m|\G\s*([a-z_]\w*)\b|igc) {
37 return ('NAME', $1);
38 }
39 elsif ($input =~ m|\G\s*([,;])|gc) {
40 return ($1, $1);
41 }

```

```

42 elsif ($input =~ m/\G\s*(.)/gc) {
43 die "Caracter invalido: $1\n";
44 }
45 else {
46 return ('', undef); # end of file
47 }
48 redo;
49 }
50 }
51 }
52
53 my $debug_level = (@ARGV)? oct(shift @ARGV): 0x1F;
54 my $parser = Inherited->new();
55 $parser->YYParse(yylex => \&scanner, yyerror => \&Error, yydebug => $debug_level);

```

En las líneas de la 15 a la 51 esta nuestro analizador léxico. La entrada se lee en una variable local cuyo valor permanece entre llamadas: hemos creado una clausura con la variable `$input` (véase la sección [?] para mas detalles sobre el uso de clausuras en Perl). Aunque la variable `$input` queda inaccesible desde fuera de la clausura, persiste entre llamadas como consecuencia de que la subrutina `scanner` la utiliza.

A continuación sigue un ejemplo de ejecución:

```

$./useinherited.pl 0
Entrada (En Unix, presione CTRL-D para terminar):
global integer x, y, z;
local float a,b;
x de clase GLOBAL, tipo INTEGER
y de clase GLOBAL, tipo INTEGER
z de clase GLOBAL, tipo INTEGER
a de clase LOCAL, tipo FLOAT
b de clase LOCAL, tipo FLOAT

```

## 8.29. Acciones en Medio de una Regla y Atributos Heredados

La estrategia utilizada en la sección 8.28 funciona si podemos predecir la posición del atributo en la pila del analizador. En el ejemplo anterior los atributos clase y tipo estaban siempre, cualquiera que fuera la derivación a derechas, en las posiciones 0 y -1. Esto no siempre es así. Consideremos la siguiente *definición dirigida por la sintaxis*:

|                         |                       |
|-------------------------|-----------------------|
| $S \rightarrow a A C$   | $\$C\{i\} = \$A\{s\}$ |
| $S \rightarrow b A B C$ | $\$C\{i\} = \$A\{s\}$ |
| $C \rightarrow c$       | $\$C\{s\} = \$C\{i\}$ |
| $A \rightarrow a$       | $\$A\{s\} = "a"$      |
| $B \rightarrow b$       | $\$B\{s\} = "b"$      |

**Ejercicio 8.29.1.** *Determine un orden correcto de evaluación de la anterior definición dirigida por la sintaxis para la entrada `b a b c`.*

C hereda el atributo sintetizado de A. El problema es que, en la pila del analizador el atributo  $\$A\{s\}$  puede estar en la posición 0 o -1 dependiendo de si la regla por la que se derivó fue  $S \rightarrow a A C$  o bien  $S \rightarrow b A B C$ . La solución a este tipo de problemas consiste en insertar acciones intermedias de copia del atributo de manera que se garantice que el atributo de interés está siempre a una distancia

fija. Esto es, se inserta una variable sintáctica intermedia auxiliar M la cual deriva a vacío y que tiene como acción asociada una regla de copia:

|                           |                                            |
|---------------------------|--------------------------------------------|
| $S \rightarrow a A C$     | $\$C\{i\} = \$A\{s\}$                      |
| $S \rightarrow b A B M C$ | $\$M\{i\} = \$A\{s\}; \$C\{i\} = \$M\{s\}$ |
| $C \rightarrow c$         | $\$C\{s\} = \$C\{i\}$                      |
| $A \rightarrow a$         | $\$A\{s\} = "a"$                           |
| $B \rightarrow b$         | $\$B\{s\} = "b"$                           |
| $M \rightarrow \epsilon$  | $\$M\{s\} = \$M\{i\}$                      |

El nuevo esquema de traducción puede ser implantado mediante un programa **eyapp**:

```
$ cat -n Inherited2.y
1 %%
2 S : 'a' A C
3 | 'b' A B { $_[2]; } C
4 ;
5
6 C : 'c' { print "Valor: ",$_[0]->YYSemval(0),"\\n"; $_[0]->YYSemval(0) }
7 ;
8
9 A : 'a' { 'a' }
10 ;
11
12 B : 'b' { 'b' }
13 ;
14
15 %%
```

La ejecución muestra como se ha propagado el valor del atributo:

```
$./useinherited2.pl '0x04'
Entrada (En Unix, presione CTRL-D para terminar):
b a b c
Shift 2. Shift 6.
Reduce using rule 5 (A,1): Back to state 2, then state 5.
Shift 8.
Reduce 6 (B,1): Back to state 5, then state 9.
Reduce 2 (@1-3,0): Back to state 9, then state 12.
```

En este momento se esta ejecutando la acción intermedia. Lo podemos comprobar revisando el fichero **Inherited2.output** que fué generado usando la opción **-v** al llamar a **eyapp**. La regla 2 por la que se reduce es la asociada con la acción intermedia:

```
$ cat -n Inherited2.output
1 Rules:
2 -----
3 0: $start -> S $end
4 1: S -> 'a' A C
5 2: @1-3 -> /* empty */
6 3: S -> 'b' A B @1-3 C
7 4: C -> 'c'
8 5: A -> 'a'
9 6: B -> 'b'
...

```

Obsérvese la notación usada por `eyapp` para la *acción en medio de la regla*: @1-3. Continuamos con la antiderivación:

Shift 10.

Reduce 4 (C,1):

Valor: a

Back to state 12, then 13.

Reduce using rule 3 (S,5): Back to state 0, then state 1.

Shift 4.

Accept.

El método puede ser generalizado a casos en los que el atributo de interés este a diferentes distancias en diferentes reglas sin mas que introducir las correspondientes acciones intermedias de copia.



## Capítulo 9

# Análisis Semántico con Parse::Eyapp

### 9.1. Esquemas de Traducción: Conceptos

**Definición 9.1.1.** *Un esquema de traducción es una gramática independiente del contexto en la cual se han insertado fragmentos de código en las partes derechas de sus reglas de producción. Los fragmentos de código así insertados se denominan acciones semánticas. Dichos fragmentos actúan, calculan y modifican los atributos asociados con los nodos del árbol sintáctico. El orden en que se evalúan los fragmentos es el de un recorrido primero-profundo del árbol de análisis sintáctico.*

Obsérvese que, en general, para poder aplicar un esquema de traducción hay que construir el árbol sintáctico y después aplicar las acciones empotradas en las reglas en el orden de recorrido primero-profundo. Por supuesto, si la gramática es ambigua una frase podría tener dos árboles y la ejecución de las acciones para ellos podría dar lugar a diferentes resultados. Si se quiere evitar la multiplicidad de resultados (interpretaciones semánticas) es necesario precisar de que árbol sintáctico concreto se está hablando.

Por ejemplo, si en la regla  $A \rightarrow \alpha\beta$  insertamos un fragmento de código:

$$A \rightarrow \alpha\{action\}\beta$$

La acción  $\{action\}$  se ejecutará después de todas las acciones asociadas con el recorrido del subárbol de  $\alpha$  y antes que todas las acciones asociadas con el recorrido del subárbol  $\beta$ .

El siguiente esquema de traducción recibe como entrada una expresión en infijo y produce como salida su traducción a postfijo para expresiones aritmeticas con sólo restas de números:

$$\begin{array}{ll} expr \rightarrow expr_1 - NUM & \{ \$expr\{TRA\} = \$expr[1]\{TRA\}." ".$NUM\{VAL\}." - " \} \\ expr \rightarrow NUM & \{ \$expr\{TRA\} = $NUM\{VAL\} \} \end{array}$$

Las apariciones de variables sintácticas en una regla de producción se indexan como se ve en el ejemplo, para distinguir de que nodo del árbol de análisis estamos hablando. Cuando hablemos del atributo de un nodo utilizaremos una indexación tipo *hash*. Aquí VAL es un atributo de los nodos de tipo NUM denotando su valor numérico y para accederlo escribiremos \$NUM{VAL}. Análogamente \$expr{TRA} denota el atributo “traducción” de los nodos de tipo expr.

**Ejercicio 9.1.1.** *Muestre la secuencia de acciones a la que da lugar el esquema de traducción anterior para la frase 7 -5 -4.*

En este ejemplo, el cómputo del atributo \$expr{TRA} depende de los atributos en los nodos hijos, o lo que es lo mismo, depende de los atributos de los símbolos en la parte derecha de la regla de producción. Esto ocurre a menudo y motiva la siguiente definición:

**Definición 9.1.2.** *Un atributo tal que su valor en un nodo puede ser computado en términos de los atributos de los hijos del nodo se dice que es un atributo sintetizado.*

**Ejemplo 9.1.1.** *Un ejemplo de atributo heredado es el tipo de las variables en las declaraciones:*

```
decl → type { $list{T} = $type{T} } list
type → INT { $type{T} = $int }
type → STRING { $type{T} = $string }
list → ID , { $ID{T} = $list{T}; $list_1{T} = $list{T} } list1
list → ID { $ID{T} = $list{T} }
```

**Definición 9.1.3.** *Un atributo heredado es aquel cuyo valor se computa a partir de los valores de sus hermanos y de su padre.*

**Ejercicio 9.1.2.** *Escriba un esquema de traducción que convierta expresiones en infijo con los operadores  $+-*/()$  y números en expresiones en postfijo. Explique el significado de los atributos elegidos.*

## 9.2. Esquemas de Traducción con Parse::Eyapp

La distribución Parse::Eyapp debida al autor de estos apuntes (i.e. Casiano Rodriguez-Leon) permite la escritura de esquemas de traducción. El módulo no esta aún disponible en CPAN: estoy a la espera de completar las pruebas y la documentación.

El ejemplo simple que sigue ilustra como construir un esquema de traducción. El código completo puede encontrarlo en la página ???. Un ejemplo de ejecución del programa se encuentra en la página ???.

Comencemos por el principio:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> head -n79 trans_scheme_simple2.pl | cat -n
 1 #!/usr/bin/perl -w
 2 use strict;
 3 use Data::Dumper;
 4 use Parse::Eyapp;
 5 use IO::Interactive qw(is_interactive);
 6
 7 my $translationscheme = q{
 8 %{
 9 # head code is available at tree construction time
10 use Data::Dumper;
11
12 our %sym; # symbol table
13 %}
14
15 %metatree
16
17 %right '='
18 %left '- ' '+ '
19 %left '* ' '/'
20
21 %%
22 line: %name EXP
23 exp <+ ';>' /* Expressions separated by semicolons */
24 { $lhs->{n} = [map { $_->{n}} $_[1]->Children()]; }
25 ;
26
27 exp:
28 %name PLUS
```

```

29 exp.left '+' exp.right
30 { $lhs->{n} = $left->{n} + $right->{n} }
31 | %name MINUS
32 exp.left '-' exp.right
33 { $lhs->{n} = $left->{n} - $right->{n} }
34 | %name TIMES
35 exp.left '*' exp.right
36 { $lhs->{n} = $left->{n} * $right->{n} }
37 | %name DIV
38 exp.left '/' exp.right
39 { $lhs->{n} = $left->{n} / $right->{n} }
40 | %name NUM $NUM
41 { $lhs->{n} = $NUM->{attr} }
42 | '(' $exp ')' %begin { $exp }
43 | %name VAR
44 $VAR
45 { $lhs->{n} = $sym{$VAR->{attr}}->{n} }
46 | %name ASSIGN
47 $VAR '=' $exp
48 { $lhs->{n} = $sym{$VAR->{attr}}->{n} = $exp->{n} }
49
50 ;
51
52 %%
53 # tail code is available at tree construction time
54 sub _Error {
55 my($token)=$_[0]->YYCurval;
56 my($what)= $token ? "input: '$token'" : "end of input";
57
58 die "Syntax error near $what.\n";
59 }
60
61 sub _Lexer {
62 my($parser)=shift;
63
64 for ($parser->YYData->{INPUT}) {
65 $_ or return('','undef');
66
67 s/^\s*//;
68 s/^[0-9]+(?:\.[0-9]+)?// and return('NUM',$1);
69 s/^[A-Za-z][A-Za-z0-9_]*/ and return('VAR',$1);
70 s/^(.)// and return($1,$1);
71 s/^\s*//;
72 }
73 }
74
75 sub Run {
76 my($self)=shift;
77 return $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error);
78 }
79 }; # end translation scheme

```

Las líneas 7-79 no hacen otra cosa que iniciar una cadena conteniendo el esquema de traducción. También podríamos haberlo escrito en un fichero y compilarlo con `eyapp`.

## Partes de un programa eyapp

Un programa **eyapp** es similar en muchos aspectos a un programa **eyapp** / **yacc** y se divide en tres partes: la cabeza, el cuerpo y la cola. Cada una de las partes va separada de las otras por el símbolo **%%** en una línea aparte. Así, el **%%** de la línea 21 separa la cabeza del cuerpo. En la cabecera se colocan el código de inicialización, las declaraciones de terminales, las reglas de precedencia, etc. El cuerpo contiene las reglas de la gramática y las acciones asociadas. Por último, la cola de un programa **eyapp** se separa del cuerpo por otro **%%** en una línea aparte. La cola contiene las rutinas de soporte al código que aparece en las acciones así como, posiblemente, rutinas para el análisis léxico (subrutina **\_Lexer** en la línea 58) y el tratamiento de errores (subrutina **\_Error** en la línea 54). Puede encontrar una descripción de la gramática de **eyapp** usando su propia notación en la página ??.

## Esquemas de Traducción con %metatree

Por defecto **Parse::Eyapp** se comporta de manera similar a **eyapp** y **yacc** generando un analizador sintáctico LALR(1) y ejecutando las acciones empotradas según una antiderivación a derechas, esto es en un recorrido del árbol de análisis sintáctico de abajo hacia arriba y de izquierda a derecha.

Mediante la directiva **%metatree** en la línea 15 le indicamos a **Parse::Eyapp** que debe generar un esquema de traducción. En tal caso, al llamar al analizador el código empotrado no es ejecutado sino que se genera un árbol de análisis sintáctico con los correspondientes nodos de código colgando del árbol en las posiciones que les corresponden.

## Las reglas

Las reglas de producción de la gramática están entre las líneas 21 y 52. Cada entrada comienza con el nombre de la variable sintáctica y va seguida de las partes derechas de sus reglas de producción separadas por barras verticales. Opcionalmente se le puede dar un nombre a la regla de producción usando la directiva **%name**. El efecto que tiene esta directiva es bendecir - durante la fase **tree construction time** - el nodo del árbol sintáctico en una clase con nombre el argumento de **%name** (los nodos del árbol sintáctico son objetos). El código - entre llaves - puede ocupar cualquier lugar en el lado derecho.

## Nombres de los atributos

En las acciones los atributos de los nodos pueden ser accedidos usando el array mágico **@\_**. De hecho, los códigos insertados en el árbol sintáctico son convertidos en subrutinas anónimas. Así **\$\_[0]** es una referencia al nodo padre asociado con el lado izquierdo, **\$\_[1]** el asociado con el primer símbolo de la parte derecha, etc. *Esta notación posicional es confusa e induce a error: si el programador cambia la regla posteriormente insertando acciones o símbolos en la parte derecha, ¡todas las apariciones de índices en el código que se refieran a nodos a la derecha del insertado deben ser modificadas!*. Algo similar ocurre si decide suprimir una acción o un símbolo. Por ello **Parse::Eyapp** proporciona mediante la *notación punto* la posibilidad de hacer una copia automática con nombre del atributo: la notación **exp.left** indica que la variable léxica **\$left** guardará una referencia al nodo que corresponde a esta instancia de la variable sintáctica **exp**. Además **Parse::Eyapp** provee la variable léxica especial **\$lhs** donde se guarda una referencia al nodo padre. Así la regla:

```
exp.left '-' exp.right
{ $lhs->{n} = $left->{n} - $right->{n} }
```

equivale al siguiente código:

```
exp '-' exp
{
 my $lhs = shift;
 my ($left, $right) = @_[1, 3];
 $lhs->{n} = $left->{n} - $right->{n}
}
```

Si se desea usar el propio nombre de la variable sintáctica como nombre del atributo se usa la *notación dolar*. Así la notación `$exp` puede considerarse una abreviación a la notación `exp.exp`. El código:

```
$VAR '=' $exp
{ $lhs->{n} = $sym{$VAR->{attr}}->{n} = $exp->{n} }
```

equivale a este otro:

```
VAR '=' exp
{
 my $lhs = shift;
 my ($VAR, $exp) = @_[1, 3];
 $lhs->{n} = $sym{$VAR->{attr}}->{n} = $exp->{n}
}
```

## Nombres de atributos de expresiones complejas Eyapp

Consideremos la parte derecha de la regla de producción:

```
program: definition<%name PROGRAM +>.program
```

observe los siguientes detalles:

1. Es posible usar los símbolos menor/mayor para incidir en el hecho de que la especificación del tipo del nodo construido como nodo PROGRAM se le da a la lista no vacía de repeticiones de **definition**. Funciona como un paréntesis de agrupamiento y es opcional. Podríamos haber escrito:

```
program: definition %name PROGRAM +.program
```

pero es, sin duda, más confuso para el lector.

2. Si queremos darle nombre al nodo asociado con el único elemento en la parte derecha de la regla:

```
program: definition %name PROGRAM +.program
```

no podemos usar la *notación dolar*. Este código que pretende ser equivalente al anterior:

```
program:
 $definition<%name PROGRAM +>
```

produce un mensaje de error:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> eyapp Simple4
```

```
Fatal $ is allowed for identifiers only (Use dot notation instead), at line 149 at file
```

La razón del mensaje es doble:

- Primero: se está intentando dar un nombre singular **definition** a algo que es una lista y por tanto plural: debería ser algo parecido a **definitions**. La decisión del programador es inadecuada.
- Segundo: Los operadores de lista tienen mayor prioridad que el dolar: la expresión `$definition<%name PROGRAM +>` es interpretada como `$(definition<%name PROGRAM +>)`. El dolar está actuando sobre la expresión compuesta `definition<%name PROGRAM +>` que no es un identificador.

La solución: *En vez del dolar deberemos usar la notación punto para darle nombre al atributo, como se ha hecho en el ejemplo.*

## Fases de un Esquema de Traducción

La ejecución de un esquema de traducción por `Parse::Eyapp` ocurre en tres tiempos.

## Class Construction Time

En una primera parte - que denominaremos *Class Construction Time* - se analiza la gramática y se crea la clase que contendrá el analizador sintáctico. Esto se hace llamando al método de clase `new_grammar` el cual devuelve una cadena conteniendo información sobre las ambigüedades, conflictos y errores que pueda tener la gramática:

```
84 my $warnings = Parse::Eyapp->new_grammar(
85 input=>$translationscheme,
86 classname=>'main',
87 firstline => 6,
88 outputfile => 'main.pm');
89 die "$warnings\nSolve Ambiguities. See file main.output\n" if $warnings;
```

El nombre de la clase o package en el que se crea el analizador se especifica mediante el argumento `classname`.

El argumento `firstline` facilita la emisión de errores y warnings indicando la línea en que comienza la cadena que contiene el esquema de traducción.

Si se especifica el argumento `outputfile => filename` los resultados del análisis se volcarán en los ficheros `filename.pm` y `filename.output` los cuales contienen respectivamente el código del analizador e información pormenorizada sobre las tablas usadas por el analizador y los conflictos y ambigüedades encontradas durante el estudio de la gramática.

Una vez creada la clase es posible instanciar objetos del tipo analizador llamando al constructor `new` de la clase creada:

```
90 my $parser = main->new();
```

## Tree Construction Time

En una segunda parte - que denominaremos *Tree Construction Time* - se toma la entrada (usando para ello el analizador léxico y las rutinas de error proveídas por el programador) y se procede a la construcción del árbol. *Las acciones especificadas por el programador en el esquema no son ejecutadas* sino que se añaden al árbol como referencias a subrutinas (nodos de tipo CODE).

El programador puede influir en la construcción del árbol por medio de diversas directivas. De estas explicaremos tres:

### La directiva `%name class`

Como se ha dicho, la directiva `%name class` hace que el nodo asociado con la instanciación de la regla de producción se bendiga en la clase dada por la cadena `class`.

### La directiva `%begin`

La directiva `%begin { ... code ... }` usada en la línea 42 hace que el código usado como argumento `{ ... code ... }` se ejecute en *Tree Construction Time*.

```
27 exp:
28 %name PLUS
29 exp.left '+' exp.right
30 { $lhs->{n} = $left->{n} + $right->{n} }
31 | %name MINUS
..
42 | '(' $exp ')' %begin { $exp }
```

En el ejemplo la directiva `%begin { $exp }` hace que nos saltamos el nodo asociado con el paréntesis enlazando directamente la raíz del árbol referenciado por `$exp` con el padre de la regla actual. Si no se hubiera insertado esta directiva el árbol construido para la entrada `2*(3+4)` sería similar a este:

```

TIMES
|-- NUM -- TERMINAL(attr => 2)
|-- '*'
'-- E_7
 |-- '('
 |-- PLUS
 | |-- NUM -- TERMINAL(attr => 3)
 | |-- '+'
 | '--- NUM -- TERMINAL(attr => 4)
 '--- ')',

```

El efecto de la directiva `%begin { $exp }` es retornar la referencia a la expresión parentizada dando lugar al siguiente árbol:

```

TIMES
|-- NUM -- TERMINAL(attr => 2)
|-- '*'
'-- PLUS
 |-- NUM -- TERMINAL(attr => 3)
 |-- '+'
 '--- NUM -- TERMINAL(attr => 4)

```

En general, las acciones asociadas con directivas `%begin` modifican la construcción del árbol sintáctico concreto para dar lugar a un árbol de análisis sintáctico abstracto adecuado a los requerimientos de las fases posteriores.

Las acciones en Tree Construction Time insertadas mediante `%begin` se ejecutan colaborativamente con las acciones de construcción del árbol en el orden usual de los analizadores LR: según una antiderivación a derechas, esto es, en un recorrido del árbol de análisis sintáctico de abajo hacia arriba (de las hojas hacia la raíz) y de izquierda a derecha.

Las acciones en Tree Construction Time reciben como argumentos en `$_[1]`, `$_[2]`, etc. las referencias a los nodos del árbol asociadas con los elementos de la parte derecha. *En Tree Construction Time el argumento `$_[0]` es una referencia al objeto analizador sintáctico.*

La segunda fase en nuestro ejemplo ocurre en las líneas 90-92 en las que leemos la entrada y llamamos al método `Run` el cual construye el árbol:

```

90 print "Write a sequence of arithmetic expressions: " if is_interactive();
91 $parser->YYData->{INPUT} = <>;
92 my $t = $parser->Run() or die "Syntax Error analyzing input";

```

El método `Run` se limita a llamar al método `YYParse` que es quien realiza el análisis:

```

74 sub Run {
75 my($self)=shift;
76 return $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error);
77 }

```

Cuando el método `YYParse` proveído por `Parse::Eyapp` es llamado es necesario que hayan sido especificadas las correspondientes referencias a las rutinas de análisis léxico (argumento con clave `yylex`) y de tratamiento de errores (argumento con clave `yyerror`).

Después de esta fase tenemos el árbol de análisis extendido con los nodos de tipo `CODE`.

## Execution Time

En una tercera parte - que denominaremos *Execution Time* - el árbol es recorrido en orden primero-profundo y los nodos de la clase `CODE` son ejecutados. El árbol será modificado y decorado como consecuencia de las acciones y podremos examinar los resultados:

```

93 $t->translation_scheme;
94 my $treestring = Dumper($t);
95 our %sym;
96 my $symboltable = Dumper(\%sym);
97 print <<"EOR";
98 *****Tree*****
99 $treestring
100 *****Symbol table*****
101 $symboltable
102 *****Result*****
103 $t->{n}
104
105 EOR

```

El método `translation_scheme` tiene una estructura simple y tiene un código similar a este:

```

sub translation_scheme {
 my $self = shift; # root of the subtree
 my @children = $self->children();
 for (@children) {
 if (ref($_) eq 'CODE') {
 $_->($self, @children);
 }
 elsif (defined($_)) {
 translation_scheme($_);
 }
 }
}

```

Como se ve en el código de `translation_scheme` la subrutina asociada se le pasan como argumentos referencias al nodo y a los hijos del nodo.

## Los Terminales

Durante la fase de construcción del árbol sintáctico los nodos que corresponden a terminales o tokens de la gramática son -por defecto- bendecidos en la clase "``${PREFIX}` TERMINAL`". Si el programador no ha indicado lo contrario en la llamada al analizador, `$PREFIX` es la cadena vacía. (Véase el párrafo en la página 571 sobre el argumento `yyprefix` del método constructor del analizador).

Los nodos de la clase `TERMINAL` poseen al menos dos atributos `token` y `attr`. El atributo `token` indica que clase de terminal es (`NUM`, `IDENTIFIER`, etc.). El atributo `attr` nos da el valor semántico del terminal tal y como fué recibido del analizador léxico.

**Listas y Opcionales** El fragmento del esquema de traducción entre las líneas 26 y 30:

```

26 line: %name PROG
27 exp <%name EXP + ';'>
28 { @{$lhs->{t}} = map { $_->{t}} ($lhs->child(0)->Children()); }
29
30 ;

```

expresa que el lenguaje generado por el no terminal `line` esta formado por secuencias no vacías de frases generadas a partir de `exp` separadas por puntos y comas. En concreto, el analizador generado por `eyapp` transforma la regla `line: exp <%name EXP + ';'>` en:



```

line: %name EXP
 PLUS-1
;
PLUS-1: %name _PLUS_LIST
 PLUS-1 ';' exp
 |
 exp
;

```

La expresión `exp <+ ';'>` es tratada como un único elemento de la parte derecha y su atributo es un nodo de la clase `_PLUS_LIST` cuyos hijos son los elementos de la lista. Por ejemplo, para la entrada `a=2; b = 2*a` el analizador construye un árbol similar a este:

```

bless({
 'children' => [
 bless({ # _PLUS_LIST
 | 'children' => [
 | bless({ # ASSIGN a = 2
 | | 'children' => [
 | | | bless({ 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL'),
 | | | bless({ 'attr' => '=', 'token' => '=' }, 'TERMINAL'),
 | | | bless({ # NUM 2
 | | | | 'children' => [
 | | | | | bless({ 'attr' => '2', 'token' => 'NUM' }, 'TERMINAL'),
 | | | | | sub { my $lhs = $_[0]; my $NUM = $_[1]; $lhs->{n} = $NUM->{attr} }
 | | | |]
 | | | | }, 'NUM'),
 | | | | sub { my ($lhs, $exp, $VAR) = ($_[0], $_[3], $_[1]);
 | | | | | $lhs->{n} = $sym{$VAR->{attr}}->{n} = $exp->{n} }
 | | | |]
 | | | }, 'ASSIGN'),
 | | | bless({ # ASSIGN b = 2*a
 | | | | 'children' => [
 | | | | | bless({ 'attr' => 'b', 'token' => 'VAR' }, 'TERMINAL'),
 | | | | | bless({ 'attr' => '=', 'token' => '=' }, 'TERMINAL'),
 | | | | | bless({ # TIMES 2*a
 | | | | | | 'children' => [
 | | | | | | | bless({ }, 'NUM'),
 | | | | | | | bless({ 'attr' => '*', 'token' => '*' }, 'TERMINAL'),
 | | | | | | | bless({ }, 'VAR'),
 | | | | | | | sub { ... }
 | | | | | |]
 | | | | | | }, 'TIMES'),
 | | | | | | sub { ... }
 | | | | |]
 | | | | }, 'ASSIGN')
 | | |]
 | | }, '_PLUS_LIST'),
 | | sub { ... }
 |]
 }, 'EXP')

```

Observe que, por defecto, los nodos punto y coma (;) son eliminados del nodo lista de hijos del nodo `_PLUS_LIST`. Por defecto, en los nodos creados por `Parse::Eyapp` desde listas declaradas mediante operadores de brackets (por ejemplo `St <+ ';'>` o `ID <* ', '>`) se elimina el separador si este fué definido mediante una cadena (uso de apostrofes).

Diremos que un terminal es un *terminal sintáctico* o *syntax token* si fue definido en el programa `eyapp` mediante una cadena delimitada por apóstrofes.

Si queremos cambiar el estatus de un *syntax token*, por ejemplo si queremos que el separador `';` del ejemplo forme parte de la lista deberemos añadir a la cabecera la declaración `%semantic token ';'.`

Si creamos una nueva versión de nuestro programa `trans_scheme_simple3.pl` añadiendo esta declaración:

```
%semantic token ';'
%right '='
....
```

Las listas contendrán los puntos y comas. En tal caso, la línea 24 dará lugar a un error<sup>1</sup> ya que los nodos punto y coma carecen del atributo `n`:

```
24 { $lhs->{n} = [map { $_->{n}} $_[1]->Children()]; }
```

En efecto:

```
nereida:~/src/perl/YappWithDefaultAction/examples> trans_scheme_simple3.pl > salida
a=2*3; b = a+1; c = a-b
Use of uninitialized value in join or string at trans_scheme_simple3.pl line 99, <> line 1.
Use of uninitialized value in join or string at trans_scheme_simple3.pl line 99, <> line 1.
```

Al usar la declaración `%semantic token` la nueva estructura del árbol es:

```
bless({
 'n' => [6, undef, 7, undef, -1],
 'children' => [
 bless({
 'children' => [
 bless({ 'n' => 6,] }, 'ASSIGN'),
 bless({ 'children' => [], 'attr' => ';', 'token' => ';' }, 'TERMINAL'),
 bless({ 'n' => 7, }, 'ASSIGN'),
 bless({ 'children' => [], 'attr' => ';', 'token' => ';' }, 'TERMINAL'),
 bless({ 'n' => -1, }, 'ASSIGN')
]
 }, '_PLUS_LIST'),
 sub { "DUMMY" }
]
}, 'EXP')
```

`Parse::Eyapp` extiende `Parse::Yapp` con listas vacías y no vacías usando los operadores `*`, `+`:

La regla:

$$A : B C * 'd'$$

es equivalente a:

---

<sup>1</sup> El código de la línea 24

```
24 { $lhs->{n} = [map { $_->{n}} $_[1]->Children()]; }
```

funciona como sigue: obtiene la lista de hijos verdaderos (aquellos que no son referencias a subrutinas) mediante el método `Children`. La aplicación posterior de `map` crea una lista con los atributos `n` de esos nodos. Por último la inclusión entre corchetes crea una lista anónima con dichos números. La referencia resultante es asignada al atributo `n` del lado izquierdo

```

A : B L 'd'
L : /* vacío */
 | P
P : C P
 | C

```

Es posible especificar la presencia de símbolos opcionales usando ?.

Observe que el operador de concatenación tiene menor prioridad que los operadores de listas, esto es la expresión  $AB^*$  es interpretada como  $A(B^*)$ .

Una secuencia de símbolos en la parte derecha de una regla de producción puede ser agrupada mediante el uso de paréntesis. Al agrupar una secuencia se crea una variable sintáctica intermedia que produce dicha secuencia. Por ejemplo, la regla:

```
A : B (C { dosomething(@_) })? D
```

es equivalente a:

```

A : B T1 D
T1 : /* vacío */
 | T2
T2 : C { dosomething(@_) }

```

## Ambigüedades

Hay numerosas ambigüedades en la gramática asociada con el esquema de traducción presentado en la página ???. Las ambigüedades se resuelven exactamente igual que en `yacc` usando directivas en la *cabecera* o primera parte que indiquen como resolverlas.

Entre las ambigüedades presentes en la gramática del ejemplo estan las siguientes:

- ¿Como debo interpretar la expresión  $e - e - e$ ? ¿Como  $(e - e) - e$ ? ¿o bien  $e - (e - e)$ ? La respuesta la da la asignación de asociatividad a los operadores que hicimos en la cabecera. Al declarar como asociativo a izquierdas al terminal `-` hemos resuelto este tipo de ambigüedad. Lo que estamos haciendo es indicarle al analizador que a la hora de elegir entre los árboles abstractos elija siempre el árbol que se hunde a izquierdas.

- ¿Como debo interpretar la expresión  $e - e * e$ ? ¿Como  $(e - e) * e$ ? ¿o bien  $e - (e * e)$ ? En `eyapp` los terminales declarados mediante directivas `%left`, `%right` y `%nonassoc` tienen asociada una prioridad. *Esa prioridad es mayor cuanto mas abajo en el texto está la línea de su declaración.* Un terminal puede ser también declarado con la directiva `%token` en cuyo caso no se le asocia prioridad.

Al declarar que el operador `*` tiene mayor prioridad que el operador `-` estamos resolviendo esta otra fuente de ambigüedad. Esto es así pues el operador `*` fué declarado después que el operador `-`. Le indicamos así al analizador que construya el árbol asociado con la interpretación  $e - (e * e)$ .

*En `eyapp` La prioridad asociada con una regla de producción es la del último terminal que aparece en dicha regla.*

Una regla de producción puede ir seguida de una directiva `%prec` la cual le da una prioridad explícita. Esto puede ser de gran ayuda en ciertos casos de ambigüedad. Por ejemplo, si quisieramos introducir el uso el menos unario en la gramática surgiría una ambigüedad:

```

39 | %name UMINUS
40 '-' $exp %prec NEG
41 { $lhs->{n} = -$exp->{n} }

```

¿Cuál es la ambigüedad que surge con esta regla? Una de las ambigüedades de esta regla está relacionada con el doble significado del menos como operador unario y binario: hay frases como  $-e-e$  que tiene dos posibles interpretaciones: Podemos verla como  $(-e)-e$  o bien como  $-(e-e)$ . Hay dos árboles posibles. El analizador, cuando esté analizando la entrada  $-e-e$  y vea el segundo  $-$  deberá escoger uno de los dos árboles. ¿Cuál?

El conflicto puede verse como una “lucha” entre la regla `exp: '-' exp` la cual interpreta la frase como  $(-e)-e$  y la segunda aparición del terminal  $-$  el cual “quiere entrar” para que gane la regla `exp: exp '-' exp` y dar lugar a la interpretación  $-(e-e)$ .

En este caso, si atendemos a la norma enunciada de que la prioridad asociada con una regla de producción es la del último terminal que aparece en dicha regla, las dos reglas  $E \rightarrow -E$  y  $E \rightarrow E - E$  tienen la prioridad del terminal  $-$ .

Lo que hace la declaración `%prec NEG` de la línea 40 es modificar la prioridad de la regla  $E \rightarrow -E$  para que tenga la del terminal `NEG`. El terminal `NEG` lo declaramos en la cabecera del programa, dándole la prioridad adecuada:

```
17 %right '='
18 %left '-' '+'
19 %left '*' '/'
20 $right NEG
21 %%
```

**Ejercicio 9.2.1.** ¿Cómo se hubiera interpretado la expresión  $-e-e$  si no se hubiese introducido el terminal de desempate `NEG`?

**Ejercicio 9.2.2.** ¿Cómo se interpretará una expresión como  $e = e - e$ ?

- Tenga en cuenta que algunas funcionalidades proveídas por `Parse::Eyapp` (por ejemplo las listas) suponen la inserción de reglas dentro de la gramática de entrada y por tanto pueden dar lugar a ambigüedades.
- Una fuente de ambigüedad puede aparecer como consecuencia de la aparición de acciones semánticas. Cuando se inserta una acción  $\{action_1\}$  en medio de una regla  $A \rightarrow \alpha\beta$ :

$$A \rightarrow \alpha \{action_1\} \beta \{action_2\}$$

`eyapp` crea una variable sintáctica temporal  $T$  e introduce una nueva regla:

1.  $A \rightarrow \alpha T \beta \{action_2\}$
2.  $T \rightarrow \epsilon \{action_1\}$

esta modificación puede dar lugar a conflictos.

## La rutina de Tratamiento de Errores

Recuerde que la fase de análisis sintáctico y léxico de la entrada ocurre en Tree Construction Time. En consecuencia el primer argumento que recibe el método `_Error` (líneas 53-59) cuando es llamado por el analizador sintáctico es una referencia al objeto analizador sintáctico. Dicho objeto dispone de un conjunto de métodos, muchos de los cuales ya existían en `Parse::Yapp`. Entre estos últimos se encuentra el método `YYCurval` que es llamado en la línea 55 y que devuelve el terminal/token que estaba siendo analizado en el momento en el que se produjo el error. Si dicho token no está definido es que hemos alcanzado el final del fichero (línea 56).

```

54 sub _Error {
55 my($token)=$_[0]->YYCurval;
56 my($what)= $token ? "input: '$token'" : "end of input";
57
58 die "Syntax error near $what.\n";
59 }

```

Otros métodos que pueden ser de ayuda en el diagnóstico de errores son `YYCurval` que devuelve el atributo del token actual y `YYExpect` que devuelve una lista con los terminales esperados en el momento en el que se produjo el error.

## El Analizador Léxico

El analizador léxico esta tomado de los ejemplos que acompañan a `Parse::Yapp`. Se supone que la entrada se ha dejado dentro del objeto analizador en `$parser->YYData->{INPUT}`. Recuerde que el análisis léxico de la entrada ocurre en *Tree Construction Time*. En consecuencia el primer argumento que recibe `_Lexer` cuando es llamado por el analizador sintáctico es la referencia al objeto analizador sintáctico. De ahí que lo primero que se hace, en la línea 59, sea crear en `$parser` una variable léxica que referencia dicho objeto.

```

58 sub _Lexer {
59 my($parser)=shift;
60
61 $parser->YYData->{INPUT}
62 or return('',undef);
63
64 $parser->YYData->{INPUT}=~s/^\s*//;
65
66 for ($parser->YYData->{INPUT}) {
67 s/^\([0-9]+(?:\.[0-9]+)?\)// and return('NUM',$1);
68 s/^\([A-Za-z][A-Za-z0-9_]*\)// and return('VAR',$1);
69 s/^\(.\)// and return($1,$1);
70 s/^\s*//;
71 }
72 }

```

**Ejercicio 9.2.3.** *¿Cuántos elementos tiene la lista sobre la que se hace el bucle `for` de la línea 66?*

Obsérvese el *falso bucle for* en la línea 66. Es un truco que constituye una de esas frases hechas o *idioms* que aunque la primera vez resultan extrañas, a fuerza de verlas repetidas se convierten en familiares.

El bucle de hecho se ejecutará una sola vez en cada llamada a `_Lexer`. El objetivo es evitar las costosas indirecciones a las que obliga almacenar la entrada en `$parser->YYData->{INPUT}`. Para ello se aprovecha la capacidad del bucle `for` sin índice de crear en `$_` un alias del elemento visitado en la iteración.

## Ejercicio 9.2.4.

1. *¿Puedo cambiar el binding de la línea 64 por uno sobre `$_` dentro del falso `for`? ¿Ganamos algo con ello?*
2. *¿Puedo cambiar la comprobación en las líneas 61-62 por una primera línea dentro del falso `for` que diga `$_ or return('',undef);`?*
3. *¿Cuántas veces se ejecuta el falso bucle si `$parser->YYData->{INPUT}` contiene la cadena vacía?*
4. *¿Que ocurrirá en las líneas 61-62 si `$parser->YYData->{INPUT}` contiene solamente la cadena `'0'`?*

## Acciones por Defecto

En el ejemplo anterior la acción es asociada con los nodos PLUS, MINUS, TIMES y DIV es similar. Parse::Eyapp proporciona una directiva % defaultaction la cual permite especificar la acción por defecto. Esta acción es asociada con las reglas que no tienen una acción asociada explícita. El siguiente ejemplo muestra su uso:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> cat -n trans_scheme_default_action.pl
 1 #!/usr/bin/perl
 2 use strict;
 3 use warnings;
 4 use Data::Dumper;
 5 use Parse::Eyapp;
 6 use IO::Interactive qw(interactive);
 7
 8 my $translationscheme = q{
 9 %{
10 # head code is available at tree construction time
11 use Data::Dumper;
12
13 our %sym; # symbol table
14 %}
15
16 %defaultaction { $lhs->{n} = eval " $left->{n} $_[2]->{attr} $right->{n} " }
17
18 %metatree
19
20 %right '='
21 %left '- ' '+'
22 %left '* ' '/'
23
24 %%
25 line: %name EXP
26 exp <+ ';'> /* Expressions separated by semicolons */
27 { $lhs->{n} = $_[1]->Last_child->{n} }
28 ;
29
30 exp:
31 %name PLUS
32 exp.left '+' exp.right
33 | %name MINUS
34 exp.left '-' exp.right
35 | %name TIMES
36 exp.left '*' exp.right
37 | %name DIV
38 exp.left '/' exp.right
39 | %name NUM $NUM
40 { $lhs->{n} = $NUM->{attr} }
41 | '(' $exp ')' %begin { $exp }
42 | %name VAR
43 $VAR
44 { $lhs->{n} = $sym{$VAR->{attr}}->{n} }
45 | %name ASSIGN
46 $VAR '=' $exp
47 { $lhs->{n} = $sym{$VAR->{attr}}->{n} = $exp->{n} }
```

```

48
49 ;
50
51 %%
52 sub Error {
53 die "Syntax error near ".$_[0]->YYCurval?$_[0]->YYCurval:"end of file")."\n";
54 }
55
56 sub Lexer {
57 my($parser)=shift;
58
59 for ($parser->YYData->{INPUT}) {
60 s/^\s*//;
61 $_ eq '' and return('',undef);
62 s/^[0-9]+(?:\.[0-9]+)?// and return('NUM',$1);
63 s/^[A-Za-z][A-Za-z0-9_]*/ and return('VAR',$1);
64 s/^(.)// and return($1,$1);
65 }
66 }
67 }; # end translation scheme
68
69 $Data::Dumper::Indent = 1;
70 $Data::Dumper::Terse = 1;
71 $Data::Dumper::Deepcopy = 1;
72 my $warnings = Parse::Eyapp->new_grammar(
73 input=>$translationscheme,
74 classname=>'Calc',
75 firstline => 6,
76 outputfile => 'Calc.pm');
77 die "$warnings\nSolve Ambiguities. See file main.output\n" if $warnings;
78 my $parser = Calc->new();
79 print {interactive} "Write a sequence of arithmetic expressions: ";
80 $parser->YYData->{INPUT} = <>;
81 my $t = $parser->YYParse(yylex => \&Calc::Lexer, yyerror => \&Calc::Error);
82 $t->translation_scheme;
83 my $treestring = Dumper($t);
84 my $symboltable;
85 {
86 no warnings;
87 $symboltable = Dumper(\%Calc::sym);
88 }
89 print <<"EOR";
90 *****Tree*****
91 $treestring
92 *****Symbol table*****
93 $symboltable
94 *****Result*****
95 $t->{n}
96
97 EOR

```

El método `Last_child` usado en la línea 27 devuelve una referencia al último hijo no código del nodo. Al ser `$_[1]` un nodo de tipo `'PLUS_LIST'` queda garantizado que el último hijo no es una referencia a una subrutina así que podría haberse usado el método `last_child` el cual devuelve el

último hijo del nodo, sea este código o no.

La línea 86 tiene por efecto desactivar los avisos. De otra manera se produciría un `warning` con respecto al uso único de la variable `%Calc::sym`:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> echo "a=2*3; b=a+1" | trans_scheme_default_action.pl
Name "Calc::sym" used only once: possible typo at trans_scheme_default_action.pl line 85.
*****Tree*****
bless({
 'n' => 7,
 'children' => [

]
}, 'EXP')

*****Symbol table*****
{
 'a' => { 'n' => 6 },
 'b' => { 'n' => 7 }
}

*****Result*****
7
```

### Un Esquema de Traducción para las Declaraciones de Variables

En el siguiente ejemplo se muestra la implementación en `eyapp` del ejemplo 4.7.1. El código es similar salvo por la presencia de flechas de referenciado:

```
nereida:~/src/perl/YappWithDefaultAction/examples> cat -n trans_scheme_simple_decls2.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Data::Dumper;
4 use Parse::Eyapp;
5 our %s; # symbol table
6
7 my $ts = q{
8 %token FLOAT INTEGER NAME
9
10 %{
11 our %s;
12 %}
13
14 %metatree
15
16 %%
17 D1: D <* ';>'>
18 ;
19
20 D : $T { $L->{t} = $T->{t} } $L
21 ;
22
23 T : FLOAT { $lhs->{t} = "FLOAT" }
24 | INTEGER { $lhs->{t} = "INTEGER" }
25 ;
26
```



```

27 L : $NAME
28 { $NAME->{t} = $lhs->{t}; $s{$NAME->{attr}} = $NAME }
29 | $NAME { $NAME->{t} = $lhs->{t}; $L->{t} = $lhs->{t} } ', ' $L
30 { $s{$NAME->{attr}} = $NAME }
31 ;
32 %%
33 };
34
35 sub Error { die "Error sintáctico\n"; }
36
37 { # Closure of $input, %reserved_words and $validchars
38 my $input = "";
39 my %reserved_words = ();
40 my $validchars = "";
41
42 sub parametrize__scanner {
43 $input = shift;
44 %reserved_words = %{shift()};
45 $validchars = shift;
46 }
47
48 sub scanner {
49 $input =~ m{\G\s+}gc; # skip whites
50 if ($input =~ m{\G([a-zA-Z]\w*)\b}gc) {
51 my $w = uc($1); # upper case the word
52 return ($w, $w) if exists $reserved_words{$w};
53 return ('NAME', $1); # not a reserved word
54 }
55 return ($1, $1) if ($input =~ m/\G([$validchars])/gc);
56 die "Caracter invalido: $1\n" if ($input =~ m/\G(\S)/gc);
57 return ('', undef); # end of file
58 }
59 } # end closure
60
61 Parse::Eyapp->new_grammar(input=>$ts, classname=>'main', outputfile=>'Types.pm');
62 my $parser = main->new(yylex => \&scanner, yyerror => \&Error); # Create the parser
63
64 parametrize__scanner(
65 "float x,y;\ninteger a,b\n",
66 { INTEGER => 'INTEGER', FLOAT => 'FLOAT'},
67 ",;"
68);
69
70 my $t = $parser->YYParse() or die "Syntax Error analyzing input";
71
72 $t->translation_scheme;
73
74 $Data::Dumper::Indent = 1;
75 $Data::Dumper::Terse = 1;
76 $Data::Dumper::Deepcopy = 1;
77 $Data::Dumper::Deparse = 1;
78 print Dumper($t);
79 print Dumper(\%s);

```

Al ejecutarlo con la entrada "float x,y;\ninteger a,b\n" los contenidos finales del arbol son:

Los contenidos de la tabla de símbolos **%s** quedan como sigue:

## Prefijos

*Si se desean evitar posibles colisiones con clases existentes es posible prefijar todos los nombres de las clases con un prefijo dado usando el parámetro `yyprefix` en la llamada al constructor del analizador:*

```
my $warnings = Parse::Eyapp->new_grammar(
 input=>$translationscheme,
 classname=>'main',
 firstline => 6,
 outputfile => 'main.pm');
die "$warnings\nSolve Ambiguities. See file main.output\n" if $warnings;

Prefix all the classes with 'Calc::'
my $parser = main->new(yyprefix => 'Calc::');
```

El resultado de esta llamada a `new` es que las clases de los nodos quedan prefijadas con `Calc::`. Por ejemplo el árbol creado para la frase `a=1` será:

```
bless({ 'children' => [
 bless({ 'children' => [
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'Calc::TERMINAL'),
 bless({ 'children' => [], 'attr' => '=', 'token' => '=' }, 'Calc::TERMINAL'),
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => '1', 'token' => 'NUM' }, 'Calc::TERMINAL')
]
 }, 'Calc::NUM'),
]
 }, 'Calc::ASSIGN')
]
}, 'Calc::_PLUS_LIST'),
], 'Calc::EXP')
```

## Modo Standalone

Es mas eficiente aislar el código del esquema de traducción en un fichero (se asume por defecto el tipo `.eypp`) y generar el módulo que contiene el código del analizador usando el guión `eyapp`.

El siguiente ejemplo muestra un ejemplo de compilación separada. El esquema de traducción convierte expresiones en infijo a postfijo. De un lado tenemos el fichero `TSPostfix2.eypp` conteniendo el esquema:

```
nereida:~/src/perl/YappWithDefaultAction/examples> cat -n TSPostfix2.eypp
1 # File TSPostfix2.eypp
2 %right '='
3 %left '-' '+'
4 %left '*' '/'
5 %left NEG
6
7 %{
8 use Data::Dumper;
9 $Data::Dumper::Indent = 1;
10 $Data::Dumper::Deepcopy = 1;
11 #$Data::Dumper::Deparse = 1;
12 use IO::Interactive qw(interactive);
13 %}
```

```

14
15 %metatree
16
17 %defaultaction {
18 if (@_==4) { # binary operations: 4 = lhs, left, operand, right
19 $lhs->{t} = "$_[1]->{t} $_[3]->{t} $_[2]->{attr}";
20 return
21 }
22 die "Fatal Error. Unexpected input\n".Dumper(@_);
23 }
24
25 %%
26 line: %name PROG
27 exp <%name EXP + ';'>
28 { @{$lhs->{t}} = map { $_->{t}} ($lhs->child(0)->Children()); }
29
30 ;
31
32 exp: NUM { $lhs->{t} = $_[1]->{attr}; }
33 | VAR { $lhs->{t} = $_[1]->{attr}; }
34 | VAR '=' exp { $lhs->{t} = "$_[1]->{attr} $_[3]->{t} =" }
35 | exp '+' exp
36 | exp '-' exp
37 | exp '*' exp
38 | exp '/' exp
39 | '-' exp %prec NEG { $_[0]->{t} = "$_[2]->{t} NEG" }
40 | '(' exp ')' %begin { $_[2] }
41 ;
42
43 %%
44
45 sub _Error {
46 my($token)=$_[0]->YYCurval;
47
48 my($what)= $token ? "input: '$token'" : "end of input";
49 die "Syntax error near $what.\n";
50 }
51
52 my $x; # Used for input
53
54 sub _Lexer {
55 my($parser)=shift;
56
57 $x =~ s/^\s+//;
58 return('','undef') if $x eq '';
59
60
61 $x =~ s/^\([0-9]+(?:\.[0-9]+)?\)// and return('NUM',$1);
62 $x =~ s/^\([A-Za-z][A-Za-z0-9_]*\)// and return('VAR',$1);
63 $x =~ s/^\(.\)//s and return($1,$1);
64 }
65
66 sub Run {

```

```

67 my($self)=shift;
68 $x = <>;
69 my $tree = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
70 #yydebug => 0xFF
71);
72
73 print Dumper($tree);
74 $tree->translation_scheme();
75 print Dumper($tree);
76 {
77 local $" = ",";
78 print "Translation:\n@{$tree->{t}}\n";
79 }
80 }

```

Observe el uso del método `Children` - con `C` mayúscula - en la línea 28: *devuelve los hijos del nodo no incluyendo las referencias a subrutinas*. A diferencia de su homólogo con `c` minúscula `children` el cual devuelve todos los hijos del nodo, incluyendo las referencias al código empotrado.

Compilamos con `eyapp`:

```

nereida:~/src/perl/YappWithDefaultAction/examples> eyapp TSPostfix2
nereida:~/src/perl/YappWithDefaultAction/examples> ls -ltr | tail -2
-rw-r----- 1 pl users 1781 2006-10-30 13:08 TSPostfix2.eyp
-rw-r--r-- 1 pl users 7611 2006-10-30 13:12 TSPostfix2.pm

```

De otro lado tenemos el programa cliente el cual se limita a cargar el módulo y llamar al método `Run`:

```

nereida:~/src/perl/YappWithDefaultAction/examples> cat -n usetspostfix2.pl
 1 #!/usr/bin/perl -w
 2 use strict;
 3 use TSPostfix2;
 4
 5 my $parser = new TSPostfix2();
 6 $parser->Run;

```

Al ejecutar el programa se produce una salida similar a esta (la salida ha sido editada para darle mayor claridad):

```

nereida:~/src/perl/YappWithDefaultAction/examples> usetspostfix2.pl
1 a=2
2 ...
3 $VAR1 = bless({ 'children' => [
4 bless({ 'children' => [
5 | bless({ 'children' => [
6 | | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL'),
7 | | bless({ 'children' => [], 'attr' => '=', 'token' => '=' }, 'TERMINAL'),
8 | | bless({ 'children' => [
9 | | bless({ 'children' => [], 'attr' => '2', 'token' => 'NUM' }, 'TERMINAL')
10 | | sub { "DUMMY" }
11 | |],
12 | | 't' => '2'
13 | | }, 'exp_4'),
14 | | sub { "DUMMY" }
15 | |],

```

```

16 | | 't' => 'a 2 ='
17 | }, 'exp_6')
18 |]
19 }, 'EXP'),
20 sub { "DUMMY" }
21],
22 't' => ['a 2 =']
23 }, 'PROG');
24 Translation:
25 a 2 =

```

Como puede verse en la salida, cuando no se especifica el nombre del nodo asociado con la regla de producción se genera un *nombre por defecto* que consiste en la concatenación del nombre de la variable sintáctica en el lado izquierdo y el número de orden de la regla de producción. Así el nodo descrito en las líneas 5-17 tiene por nombre `exp_6` indicando que corresponde a la sexta regla de producción de `exp`.

El nodo en las líneas 4-19 tiene por nombre `EXP`. Su nombre le fué dado mediante la directiva `%name` en la línea

```

27 exp <%name EXP + ';' >

```

Es posible insertar una directiva `%name` en una lista usando esta sintaxis.

Nótese también que si hubieramos usado la opción `Data::Dumper::Deparse` (línea 11) podríamos hacer que `Data::Dumper` nos informe no sólo de la presencia de código sino que nos muestre el código fuente que ocupa esa posición.

## 9.3. Definición Dirigida por la Sintaxis y Gramáticas Atribuidas

### Definición Dirigida por la Sintaxis

**Definición 9.3.1.** Una definición dirigida por la sintaxis es un pariente cercano de los esquemas de traducción. En una definición dirigida por la sintaxis una gramática  $G = (V, \Sigma, P, S)$  se aumenta con nuevas características:

- A cada símbolo  $S \in V \cup \Sigma$  de la gramática se le asocian cero o mas atributos. Un atributo queda caracterizado por un identificador o nombre y un tipo o clase. A este nivel son atributos formales, como los parámetros formales, en el sentido de que su realización se produce cuando el nodo del árbol es creado.
- A cada regla de producción  $A \rightarrow X_1 X_2 \dots X_n \in P$  se le asocian un conjunto de reglas de evaluación de los atributos o reglas semánticas que indican que el atributo en la parte izquierda de la regla semántica depende de los atributos que aparecen en la parte derecha de la regla semántica. El atributo que aparece en la parte izquierda de la regla semántica puede estar asociado con un símbolo en la parte derecha de la regla de producción.
- Los atributos de cada símbolo de la gramática  $X \in V \cup \Sigma$  se dividen en dos grupos disjuntos: atributos sintetizados y atributos heredados. Un atributo de  $X$  es un atributo heredado si depende de atributos de su padre y hermanos en el árbol. Un atributo sintetizado es aquél tal que el valor del atributo depende de los valores de los atributos de los hijos, es decir en tal caso  $X$  ha de ser una variable sintáctica y los atributos en la parte derecha de la regla semántica deben ser atributos de símbolos en la parte derecha de la regla de producción asociada.
- Los atributos predefinidos se denominán atributos intrínsecos. Ejemplos de atributos intrínsecos son los atributos sintetizados de los terminales, los cuáles se han computado durante la fase de análisis léxico. También son atributos intrínsecos los atributos heredados del símbolo de arranque, los cuales son pasados como parámetros al comienzo de la computación.

La diferencia principal con un esquema de traducción está en que *no se especifica el orden de ejecución de las reglas semánticas*. Se asume que, bien de forma manual o automática, se resolverán las dependencias existentes entre los atributos determinadas por la aplicación de las reglas semánticas, de manera que serán evaluados primero aquellos atributos que no dependen de ningún otro, después los que dependen de estos, etc. siguiendo un esquema de ejecución que viene guiado por las dependencias existentes entre los datos.

### Evaluación de una Definición Dirigida por la Sintaxis

Aunque hay muchas formas de realizar un evaluador de una definición dirigida por la sintaxis, conceptualmente, tal evaluador debe:

1. Construir el árbol de análisis sintáctico para la gramática y la entrada dadas.
2. Analizar las reglas semánticas para determinar los atributos, su clase y las dependencias entre los mismos.
3. Construir explícita o implícitamente el *grafo de dependencias* de los atributos, el cual tiene un nodo para cada ocurrencia de un atributo en el árbol de análisis sintáctico etiquetado con dicho atributo. El grafo tiene una arista entre dos nodos si existe una dependencia entre los dos atributos a través de alguna regla semántica.
4. Supuesto que el grafo de dependencias determina un *orden parcial* (esto es cumple las propiedades reflexiva, antisimétrica y transitiva) construir un *orden topológico* compatible con el orden parcial.
5. Evaluar las reglas semánticas de acuerdo con el orden topológico.

### Gramáticas Atribuídas

**Definición 9.3.2.** Una definición dirigida por la sintaxis en la que las reglas semánticas no tienen efectos laterales se denomina una gramática atribuída.

### Gramáticas L-Atribuídas

**Definición 9.3.3.** Si la definición dirigida por la sintaxis puede ser realizada mediante un esquema de traducción se dice que es L-atribuída.

Para que una definición dirigida por la sintaxis sea L-atribuída deben cumplirse que cualquiera que sea la regla de producción  $A \rightarrow X_1 \dots X_n$ , los atributos heredados de  $X_j$  pueden depender únicamente de:

1. Los atributos de los símbolos a la izquierda de  $X_j$
2. Los atributos heredados de  $A$

Si la gramática es LL(1), resulta fácil realizarla en un analizador descendente recursivo predictivo. La implantación de atributos del cálculo de atributos heredados en un programa **eyapp** requiere un poco más de habilidad. Véase la sección 8.28.

### Gramáticas S-Atribuídas

**Definición 9.3.4.** Si la definición dirigida por la sintaxis sólo utiliza atributos sintetizados se denomina S-atribuída.

Una definición S-atribuída puede ser fácilmente trasladada a un programa **eyapp**.

Es posible realizar tanto gramáticas L-atribuídas como S-atribuídas usando un esquema de traducción **Parse::Eyapp** generado con la directiva **%metatree** (véase la sección 9.2) como combinaciones de ambas.

## 9.4. El módulo `Language::AttributeGrammar`

El módulo `Language::AttributeGrammar` de Luke Palmer permite experimentar en Perl con una variante de las gramáticas atribuidas.

### Ejemplo Simple

`Language::AttributeGrammar` provee un constructor `new` que recibe la cadena describiendo la gramática atribuida y un método `apply` que recibe como argumentos el árbol AST y el atributo a evaluar. `Language::AttributeGrammar` toma como entrada un AST en vez de una descripción textual (línea 36).

```
nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples> cat -n atg.pl
 1 #!/usr/bin/perl -w
 2 use strict;
 3 use Language::AttributeGrammar;
 4 use Data::Dumper;
 5
 6 my $grammar = new Language::AttributeGrammar <<'EOG';
 7
 8 # find the minimum from the leaves up
 9 Leaf: $/.min = { $<value> }
10 Branch: $/.min = {
11 $<left>.min <= $<right>.min ? $<left>.min : $<right>.min;
12 }
13
14 # propagate the global minimum downward
15 ROOT: $/.gmin = { $/.min }
16 Branch: $<left>.gmin = { $/.gmin }
17 Branch: $<right>.gmin = { $/.gmin }
18
19 # reconstruct the minimized result
20 Leaf: $/.result = { bless { value => $/.gmin } => 'Leaf' }
21 Branch: $/.result = { bless { left => $<left>.result,
22 right => $<right>.result } => 'Branch' }
23
24 EOG
25
26 sub Leaf { bless { value => $_[0] } => 'Leaf' }
27 sub Branch { bless { left => $_[0], right => $_[1] } => 'Branch' }
28
29 my $tree = Branch(
30 Branch(Leaf(2), Leaf(3)),
31 Branch(Leaf(1), Branch(Leaf(5), Leaf(9))));
32 my $result = Branch(
33 Branch(Leaf(1), Leaf(1)),
34 Branch(Leaf(1), Branch(Leaf(1), Leaf(1))));
35
36 my $t = $grammar->apply($tree, 'result');
37
38 $Data::Dumper::Indent = 1;
39 print Dumper($t);
```

El objeto representando a la gramática atribuida es creado mediante la llamada a `Language::AttributeGrammar::new` en la línea 6. La gramática es una secuencia de reglas semánticas. Cada regla semántica tiene la forma:



```

nodetype1 : $/.attr_1 = { CÓDIGO PERL($<hijo_i>.attr_k)}
 | $<hijo_1>.attr_2 = { CÓDIGO PERL($<hijo_i>.attr_k, $/.attr_s)}
 | $<hijo_2>.attr_3 = { CÓDIGO PERL }

nodetype2 : $/.attr_1 = { CÓDIGO PERL($<hijo_i>.attr_k)}
 | $<hijo_1>.attr_2 = { CÓDIGO PERL($<hijo_i>.attr_k, $/.attr_s)}
 | $<hijo_2>.attr_3 = { CÓDIGO PERL }


```

## Notaciones Especiales

Dentro de la especificación de la gramática es posible hacer uso de las notaciones especiales:

- `$/` se refiere al nodo que esta siendo visitado
- `$/attr` se refiere al atributo `attr` del nodo que esta siendo visitado
- `$<ident>` se refiere al hijo del nodo visitado denominado `ident`. Se asume que el nodo dispone de un método con ese nombre para obtener el hijo. La notación `$<ident>.attr` se refiere al atributo `attr` del hijo `ident` del nodo visitado.

La llamada al método `$grammar->apply($tree, 'result')` en la línea 36 dispara la computación de las reglas para el cálculo del atributo `result` sobre el árbol `$tree`.

Cuando ejecutamos el programa obtenemos la salida:

```

nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples> atg.pl
$VAR1 = bless({
 'left' => bless({
 'left' => bless({ 'value' => 1 }, 'Leaf'),
 'right' => bless({ 'value' => 1 }, 'Leaf')
 }, 'Branch'),
 'right' => bless({
 'left' => bless({ 'value' => 1 }, 'Leaf'),
 'right' => bless({
 'left' => bless({ 'value' => 1 }, 'Leaf'),
 'right' => bless({ 'value' => 1 }, 'Leaf')
 }, 'Branch')
 }, 'Branch')
}, 'Branch');

```

## Descripción del Lenguaje de Language::AttributeGrammar

Sigue una descripción en estilo `Parse::Eyapp` de la gramática aceptada por `Language::AttributeGrammar`:

```

grammar: rule *

rule: nodetype ':' (target '.' attr '=' '{' BLOQUE DE CÓDIGO PERL '}') <* '|'>

target: self | child | accesscode

```

Las variables `nodetype`, `attr`, `self`, `child` y `accesscode` viene definidas por las siguientes expresiones regulares:

```

nodetype: /(::)?\w+(::\w+)*/ # identificador Perl del tipo de nodo (PLUS, MINUS, etc.)
attr: /\w+/ # identificador del atributo
self: '$/' # Se refiere al nodo visitado
child: /\$<\w+>/ # Hijo del nodo visitado
accesscode: /'.*?'/ # Código explícito de acceso al hijo del nodo

```

## 9.5. Usando Language::AttributeGrammars con Parse::Eyapp

### La cláusula alias de %tree

El requerimiento de Language::AttributeGrammars de que los hijos de un nodo puedan ser accedidos a través de un método cuyo nombre sea el nombre del hijo no encaja con la forma oficial de acceso a los hijos que usa Parse::Eyapp. La directiva %tree de Parse::Eyapp dispone de la opción alias la cual hace que los identificadores que aparecen en los usos de la *notación punto* y de la *notación dolar* se interpreten como nombres de métodos que proveen acceso al hijo correspondiente (véase la sección 8.16)

```
nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples> cat -n Calc.pl
 1 #!/usr/bin/perl -w
 2 use strict;
 3 use Parse::Eyapp;
 4 use Data::Dumper;
 5 use Language::AttributeGrammar;
 6
 7 my $grammar = q{
.. # precedence declarations
!15 %tree bypass alias
16
17 %%
18 line: $exp { $_[1] }
19 ;
20
21 exp:
22 %name NUM
23 $NUM
24 | %name VAR
25 $VAR
26 | %name ASSIGN
27 $VAR '=' $exp
28 | %name PLUS
29 exp.left '+' exp.right
30 | %name MINUS
31 exp.left '-' exp.right
32 | %name TIMES
33 exp.left '*' exp.right
34 | %name DIV
35 exp.left '/' exp.right
!36 | %no bypass UMINUS
37 '-' $exp %prec NEG
38 | '(' $exp ')' { $_[2] } /* Let us simplify a bit the tree */
39 ;
40
41 %%
.. # as usual
78 }; # end grammar
79
80
81 $Data::Dumper::Indent = 1;
82 Parse::Eyapp->new_grammar(
83 input=>$grammar,
84 classname=>'Rule6',
```

```

85 firstline =>7,
86 outputfile => 'Calc.pm',
87);
88 my $parser = Rule6->new();
89 $parser->YYData->{INPUT} = "a = -(2*3+5-1)\n";
90 my $t = $parser->Run;
91 print "\n***** Before *****\n";
92 print Dumper($t);
93
! 94 my $attgram = new Language::AttributeGrammar <<'EOG';
! 95
! 96 # Compute the expression
! 97 NUM: $/.val = { $<attr> }
! 98 TIMES: $/.val = { $<left>.val * $<right>.val }
! 99 PLUS: $/.val = { $<left>.val + $<right>.val }
!100 MINUS: $/.val = { $<left>.val - $<right>.val }
!101 UMINUS: $/.val = { -$<exp>.val }
!102 ASSIGN: $/.val = { $<exp>.val }
!103 EOG
!104
!105 sub NUM::attr {
!106 $_[0]->{attr};
!107 }
108
!109 my $res = $attgram->apply($t, 'val');
110
111 $Data::Dumper::Indent = 1;
112 print "\n***** After *****\n";
113 print Dumper($t);
114 print Dumper($res);

```

En la línea 36 hemos usado la directiva `%no bypass` para evitar el puenteo automático del menos unario que causa la clausula `bypass` que cualifica a la directiva `%tree` de la línea 15 (repase la sección 8.15). El método `NUM::attr` en las líneas 105-107 no es necesario y podría haber sido obviado.

## Redefiniendo el acceso a los hijos

Es posible evitar el uso de la opción `alias` y forzar a `Language::AttributeGrammar` a acceder a los hijos de un nodo `Parse::Eyapp::Node` redefiniendo el método `Language::AttributeGrammar::Parser::_get_child` el cual es usado por `Language::AttributeGrammar` cuando quiere acceder a un atributo de un nodo.

```

pl@nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples$ \
cat -n CalcNumbers4.pl

1 #!/usr/bin/perl -w
2 use strict;
3 use warnings;
4 use Parse::Eyapp;
5 use Parse::Eyapp::Base qw{push_method pop_method};
6 use Data::Dumper;
7 use Language::AttributeGrammar;
8
9 my $grammar = q{
10 %right '='
11 %left '-' '+'
12 %left '*' '/'

```

```

13 %left NEG
14 %tree bypass
15
16 %%
17 line: $exp { $_[1] }
18 ;
19
20 exp:
21 %name NUM
22 $NUM
23 | %name VAR
24 $VAR
25 | %name ASSIGN
26 $var '=' $exp
27 | %name PLUS
28 exp.left '+' exp.right
29 | %name MINUS
30 exp.left '-' exp.right
31 | %name TIMES
32 exp.left '*' exp.right
33 | %name DIV
34 exp.left '/' exp.right
35 | %no bypass UMINUS
36 '-' $exp %prec NEG
37 | '(' $exp ')' { $_[2] } /* Let us simplify a bit the tree */
38 ;
39
40 var: %name VAR
41 VAR
42 ;
43 %%
44
45 ... # tail as usual
80 }; # end grammar
81
82
83 $Data::Dumper::Indent = 1;
84 Parse::Eyapp->new_grammar(input=>$grammar, classname=>'SmallCalc');
85 my $parser = SmallCalc->new();
86 $parser->YYData->{INPUT} = "a = -(2*3+ b = 5-1)\n";
87 my $t = $parser->Run;
88
89 print "\n***** Syntax Tree *****\n";
90 push_method(qw{NUM VAR}, info => sub { $_->{attr} });
91 print $t->str;
92
93 #symbol table
94 our %s;
95
96 my $attgram = new Language::AttributeGrammar <<'EOG';
97
98 # Compute the expression
99 NUM: $/.val = { $/->{attr} }

```

```

100 TIMES: $/.val = { $<0>.val * $<1>.val }
101 PLUS: $/.val = { $<0>.val + $<1>.val }
102 MINUS: $/.val = { $<0>.val - $<1>.val }
103 UMINUS: $/.val = { -$<0>.val }
104 ASSIGN: $/.val = { $<1>.val }
105 ASSIGN: $/.entry = { $::s{$<0>->{attr}} = $<1>.val; $<0>->{attr} }
106 EOG
107
108 # rewrite _get_child, save old version
109 push_method(
110 'Language::AttributeGrammar::Parser',
111 _get_child => sub {
112 my ($self, $child) = @_;
113
114 $self->child($child);
115 }
116);
117
118 my $res = $attgram->apply($t, 'entry');
119
120 # Restore old version of Language::AttributeGrammar::Parser::_get_child
121 pop_method('Language::AttributeGrammar::Parser', '_get_child');
122
123 $Data::Dumper::Indent = 1;
124 print "\n***** After call to apply*****\n";
125 push_method(qw{TIMES PLUS MINUS UMINUS ASSIGN},
126 info => sub { defined($_->{val})? $_->{val} : "" }
127);
128 print $t->str."\n"; # Tree $t isn't modified nor decorated
129
130 #The results are held in the symbol table
131 print "$_ = ${$_}\n" for keys(%s);

```

## Efectos Laterales

Cuando se ejecuta este programa el volcado de la tabla de símbolos se produce la salida:

```

pl@nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples$ \
CalcNumbers4.pl

```

```

***** Syntax Tree *****
ASSIGN(VAR[a],UMINUS(PLUS(TIMES(NUM[2],NUM[3]),ASSIGN(VAR[b],MINUS(NUM[5],NUM[1])))))
***** After call to apply*****
ASSIGN(VAR[a],UMINUS(PLUS(TIMES(NUM[2],NUM[3]),ASSIGN(VAR[b],MINUS(NUM[5],NUM[1])))))
a = -10

```

Obsérvese que el árbol original no ha cambiado después de la llamada a `apply`.

Teniendo en cuenta que la entrada era `a = -(2*3+b=5-1)` se puede sorprender que `b` no haya sido introducido en la tabla de símbolos. La razón para la ausencia de `b` en la tabla de símbolos reside en el hecho de que el atributo `entry` para el nodo raíz del árbol asociado con `b=5-1` nunca es evaluado, ya que no existe un camino de dependencias desde el atributo solicitado en la llamada a `apply` hasta el mismo. El problema se arregla creando una dependencia que de lugar a la conectividad deseada en el grafo de dependencias y provoque la ejecución del código de inserción en la tabla de símbolos como un efecto lateral.

```

nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples> \

```

```

 sed -ne '108,119p' CalcNumbers3.pl | cat -n
1 my $attgram = new Language::AttributeGrammar <<'EOG';
2
3 # Compute the expression
4 NUM: $/.val = { $/->{attr} }
5 TIMES: $/.val = { $<0>.val * $<1>.val }
6 PLUS: $/.val = { $<0>.val + $<1>.val }
7 MINUS: $/.val = { $<0>.val - $<1>.val }
8 UMINUS: $/.val = { -$<0>.val }
! 9 ASSIGN: $/.val = { $::s{$<0>->{attr}} = $<1>.val; $<1>.val }
10 EOG
11
12 my $res = $attgram->apply($t, 'val');
```

```

nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples> \
 CalcNumbers3.pl
```

```

a = -10
b = 4
```

El código de la línea 9 indica que para calcular el atributo `val` de un nodo de tipo `ASSIGN` es necesario no sólo calcular el atributo `val` del segundo hijo (`$<1>.val`) sino también ejecutar el código de inserción `$::s{$<0>->{attr}} = $<1>.val`.

## Capítulo 10

# Transformaciones Árbol con Parse::Eyapp

### 10.1. Árbol de Análisis Abstracto

Un *árbol de análisis abstracto* (denotado *AAA*, en inglés *abstract syntax tree* o *AST*) porta la misma información que el árbol de análisis sintáctico pero de forma mas condensada, eliminándose terminales y producciones que no aportan información.

#### Alfabeto con Aridad o Alfabeto Árbol

**Definición 10.1.1.** Un alfabeto con función de aridad es un par  $(\Sigma, \rho)$  donde  $\Sigma$  es un conjunto finito y  $\rho$  es una función  $\rho : \Sigma \rightarrow \mathbb{N}_0$ , denominada función de aridad. Denotamos por  $\Sigma_k = \{a \in \Sigma : \rho(a) = k\}$ .

**Lenguaje de los Arboles** Definimos el lenguaje árbol homogéneo  $B(\Sigma)$  sobre  $\Sigma$  inductivamente:

- Todos los elementos de aridad 0 están en  $B(\Sigma)$ :  $a \in \Sigma_0$  implica  $a \in B(\Sigma)$
- Si  $b_1, \dots, b_k \in B(\Sigma)$  y  $f \in \Sigma_k$  es un elemento  $k$ -ario, entonces  $f(b_1, \dots, b_k) \in B(\Sigma)$

Los elementos de  $B(\Sigma)$  se llaman árboles o términos.

**Ejemplo 10.1.1.** Sea  $\Sigma = \{A, CONS, NIL\}$  con  $\rho(A) = \rho(NIL) = 0$ ,  $\rho(CONS) = 2$ . Entonces  $B(\Sigma) = \{A, NIL, CONS(A, NIL), CONS(NIL, A), CONS(A, A), CONS(NIL, NIL), \dots\}$

**Ejemplo 10.1.2.** Una versión simplificada del alfabeto con aridad en el que estan basados los árboles construidos por el compilador de Tutu es:

$\Sigma = \{ID, NUM, LEFTVALUE, STR, PLUS, TIMES, ASSIGN, PRINT\}$

$\rho(ID) = \rho(NUM) = \rho(LEFTVALUE) = \rho(STR) = 0$

$\rho(PRINT) = 1$

$\rho(PLUS) = \rho(TIMES) = \rho(ASSIGN) = 2$ .

Observe que los elementos en  $B(\Sigma)$  no necesariamente son árboles “correctos”. Por ejemplo, el árbol  $ASSIGN(NUM, PRINT(ID))$  es un elemento de  $B(\Sigma)$ .

#### Gramática Árbol

**Definición 10.1.2.** Una gramática árbol regular es una cuadrupla  $((\Sigma, \rho), N, P, S)$ , donde:

- $(\Sigma, \rho)$  es un alfabeto con aricidad  $\rho : \Sigma \rightarrow \mathbb{N}$
- $N$  es un conjunto finito de variables sintácticas o no terminales
- $P$  es un conjunto finito de reglas de producción de la forma  $X \rightarrow s$  con  $X \in N$  y  $s \in B(\Sigma \cup N)$
- $S \in N$  es la variable o símbolo de arranque

## Lenguaje Generado por una Gramática Árbol

**Definición 10.1.3.** Dada una gramática  $(\Sigma, N, P, S)$ , se dice que un árbol  $t \in B(\Sigma \cup N)$  es del tipo  $(X_1, \dots, X_k)$  si el  $j$ -ésimo noterminal, contando desde la izquierda, que aparece en  $t$  es  $X_j \in N$ .

Si  $p = X \rightarrow s$  es una producción y  $s$  es de tipo  $(X_1, \dots, X_n)$ , diremos que la producción  $p$  es de tipo  $(X_1, \dots, X_n) \rightarrow X$ .

**Definición 10.1.4.** Consideremos un árbol  $t \in B(\Sigma \cup N)$  que sea del tipo  $(X_1, \dots, X_n)$ , esto es las variables sintácticas en el árbol leídas de izquierda a derecha son  $(X_1, \dots, X_n)$ .

- Si  $X_i \rightarrow s_i \in P$  para algún  $i$ , entonces decimos que el árbol  $t$  deriva en un paso en el árbol  $t'$  resultante de sustituir el nodo  $X_i$  por el árbol  $s_i$  y escribiremos  $t \Rightarrow t'$ . Esto es,  $t' = t\{X_i/s_i\}$
- Todo árbol deriva en cero pasos en si mismo  $t \xRightarrow{0} t$ .
- Decimos que un árbol  $t$  deriva en  $n$  pasos en el árbol  $t'$  y escribimos  $t \xRightarrow{n} t'$  si  $t$  deriva en un paso en un árbol  $t''$  el cual deriva en  $n - 1$  pasos en  $t'$ . En general, si  $t$  deriva en un cierto número de pasos en  $t'$  escribiremos  $t \xRightarrow{*} t'$ .

**Definición 10.1.5.** Se define el lenguaje árbol generado por una gramática  $G = (\Sigma, N, P, S)$  como el lenguaje  $L(G) = \{t \in B(\Sigma) : \exists S \xRightarrow{*} t\}$ .

**Ejemplo 10.1.3.** Sea  $G = (\Sigma, V, P, S)$  con  $\Sigma = \{A, CONS, NIL\}$  y  $\rho(A) = \rho(NIL) = 0$ ,  $\rho(CONS) = 2$  y sea  $V = \{exp, list\}$ . El conjunto de producciones  $P$  es:

$$P_1 = \{list \rightarrow NIL, list \rightarrow CONS(exp, list), exp \rightarrow A\}$$

La producción  $list \rightarrow CONS(exp, list)$  es del tipo  $(exp, list) \rightarrow list$ .

El lenguaje generado por  $G$  se obtiene realizando sustituciones sucesivas (derivando) desde el símbolo de arranque hasta producir un árbol cuyos nodos estén etiquetados con elementos de  $\Sigma$ . En este ejemplo,  $L(G)$  es el conjunto de arboles de la forma:

$$L(G) = \{NIL, CONS(A, NIL), CONS(A, CONS(A, NIL)), \dots\}$$

**Ejercicio 10.1.1.** Construya una derivación para el árbol  $CONS(A, CONS(A, NIL))$ . ¿De que tipo es el árbol  $CONS(exp, CONS(A, CONS(exp, L)))$ ?

Cuando hablamos del AAA producido por un analizador sintáctico, estamos en realidad hablando de un lenguaje árbol cuya definición precisa debe hacerse a través de una gramática árbol regular. Mediante las gramáticas árbol regulares disponemos de un mecanismo para describir formalmente el lenguaje de los AAA que producirá el analizador sintáctico para las sentencias Tutu.

**Ejemplo 10.1.4.** Sea  $G = (\Sigma, V, P, S)$  con

$$\begin{aligned} \Sigma &= \{ID, NUM, LEFTVALUE, STR, PLUS, TIMES, ASSIGN, PRINT\} \\ \rho(ID) &= \rho(NUM) = \rho(LEFTVALUE) = \rho(STR) = 0 \\ \rho(PRINT) &= 1 \\ \rho(PLUS) &= \rho(TIMES) = \rho(ASSIGN) = 2 \\ V &= \{st, expr\} \end{aligned}$$

y las producciones:

$$P = \left\{ \begin{array}{ll} st & \rightarrow ASSIGN(LEFTVALUE, expr) \\ st & \rightarrow PRINT(expr) \\ expr & \rightarrow PLUS(expr, expr) \\ expr & \rightarrow TIMES(expr, expr) \\ expr & \rightarrow NUM \\ expr & \rightarrow ID \\ expr & \rightarrow STR \end{array} \right\}$$



Entonces el lenguaje  $L(G)$  contiene árboles como el siguiente:

```

ASSIGN (
 LEFTVALUE,
 PLUS (
 ID,
 TIMES (
 NUM,
 ID
)
)
)

```

El cual podría corresponderse con una sentencia como  $a = b + 4 * c$ .

El lenguaje de árboles descrito por esta gramática árbol es el lenguaje de los AAA de las sentencias de Tutu.

**Ejercicio 10.1.2.** Redefina el concepto de árbol de análisis concreto dado en la definición 8.1.7 utilizando el concepto de gramática árbol. Con mas precisión, dada una gramática  $G = (\Sigma, V, P, S)$  defina una gramática árbol  $T = (\Omega, N, R, U)$  tal que  $L(T)$  sea el lenguaje de los árboles concretos de  $G$ . Puesto que las partes derechas de las reglas de producción de  $P$  pueden ser de distinta longitud, existe un problema con la aricidad de los elementos de  $\Omega$ . Discuta posibles soluciones.

**Ejercicio 10.1.3.** ¿Cómo son los árboles sintácticos en las derivaciones árbol? Dibuje varios árboles sintácticos para las gramáticas introducidas en los ejemplos 4.9.3 y 4.9.4.

Intente dar una definición formal del concepto de árbol de análisis sintáctico asociado con una derivación en una gramática árbol

## Notación de Dewey o Coordenadas de un Árbol

**Definición 10.1.6.** La notación de Dewey es una forma de especificar los subárboles de un árbol  $t \in B(\Sigma)$ . La notación sigue el mismo esquema que la numeración de secciones en un texto: es una palabra formada por números separados por puntos. Así  $t/2.1.3$  denota al tercer hijo del primer hijo del segundo hijo del árbol  $t$ . La definición formal sería:

- $t/\epsilon = t$
- Si  $t = a(t_1, \dots, t_k)$  y  $j \in \{1 \dots k\}$  y  $n$  es una cadena de números y puntos, se define inductivamente el subárbol  $t/j.n$  como el subárbol  $n$ -ésimo del  $j$ -ésimo subárbol de  $t$ . Esto es:  $t/j.n = t_j/n$

El método `descendant` de los objetos `Parse::Eyapp::Node` descrito en la sección 8.22 puede verse como una implantación de la notación de Dewey.

**Ejercicio 10.1.4.** Sea el árbol:

```

t = ASSIGN (
 LEFTVALUE,
 PLUS (
 ID,
 TIMES (
 NUM,
 ID
)
)
)

```

Calcule los subárboles  $t/\epsilon$ ,  $t/2.2.1$ ,  $t/2.1$  y  $t/2.1.2$ .

## 10.2. Selección de Código y Gramáticas Árbol

La generación de código es la fase en la que a partir de la *Representación intermedia* o *IR* se genera una secuencia de instrucciones para la máquina objeto. Esta tarea conlleva diversas subtarefas, entre ellas destacan tres:

- La selección de instrucciones o selección de código,
- La asignación de registros y
- La planificación de las instrucciones.

El problema de la *selección de código* surge de que la mayoría de las máquinas suelen tener una gran variedad de instrucciones, habitualmente cientos y muchas instrucciones admiten mas de una decena de modos de direccionamiento. En consecuencia,

There Is More Than One Way To Do It (The Translation)

*Es posible asociar una gramática árbol con el juego de instrucciones de una máquina.* Las partes derechas de las reglas de producción de esta gramática vienen determinadas por el conjunto de árboles sintácticos de las instrucciones. La gramática tiene dos variables sintácticas que denotan dos tipos de recursos de la máquina: los registros representados por la variable sintáctica  $R$  y las direcciones de memoria representadas por  $M$ . Una instrucción deja su resultado en cierto lugar, normalmente un registro o memoria. La idea es que las variables sintácticas en los lados izquierdos de las reglas representan los lugares en los cuales las instrucciones dejan sus resultados.

Ademas, a cada instrucción le asociamos un coste:

| Gramática Árbol Para un Juego de Instrucciones Simple        |                   |       |
|--------------------------------------------------------------|-------------------|-------|
| Producción                                                   | Instrucción       | Coste |
| $R \rightarrow \text{NUM}$                                   | LOADC R, NUM      | 1     |
| $R \rightarrow M$                                            | LOADM R, M        | 3     |
| $M \rightarrow R$                                            | STOREM M, R       | 3     |
| $R \rightarrow \text{PLUS}(R,M)$                             | PLUSM R, M        | 3     |
| $R \rightarrow \text{PLUS}(R,R)$                             | PLUSR R, R        | 1     |
| $R \rightarrow \text{TIMES}(R,M)$                            | TIMESM R, M       | 6     |
| $R \rightarrow \text{TIMES}(R,R)$                            | TIMESR R, R       | 4     |
| $R \rightarrow \text{PLUS}(R, \text{TIMES}(\text{NUM}, R))$  | PLUSCR R, NUM, R  | 4     |
| $R \rightarrow \text{TIMES}(R, \text{TIMES}(\text{NUM}, R))$ | TIMESCR R, NUM, R | 5     |

Consideremos la IR consistente en el AST generado por el front-end del compilador para la expresión  $x+3*(7*y)$ :

$\text{PLUS}(M[x], \text{TIMES}(N[3], \text{TIMES}(N[7], M[y])))$

Construyamos una derivación a izquierdas para el árbol anterior:

| Una derivación árbol a izquierdas para $P(M, T(N, T(N, M)))$ |                                   |              |           |
|--------------------------------------------------------------|-----------------------------------|--------------|-----------|
| Derivación                                                   | Producción                        | Instrucción  | Coste     |
| $R \Rightarrow$                                              | $R \rightarrow \text{PLUS}(R,R)$  | PLUSR R, R   | 1         |
| $P(R, R) \Rightarrow$                                        | $R \rightarrow M$                 | LOADM R, M   | 3         |
| $P(M, R) \Rightarrow$                                        | $R \rightarrow \text{TIMES}(R,R)$ | TIMESR R, R  | 4         |
| $P(M, T(R, R)) \Rightarrow$                                  | $R \rightarrow \text{NUM}$        | LOADC R, NUM | 1         |
| $P(M, T(N, R)) \Rightarrow$                                  | $R \rightarrow \text{TIMES}(R,R)$ | TIMESR R, R  | 4         |
| $P(M, T(N, T(R, R))) \Rightarrow$                            | $R \rightarrow \text{NUM}$        | LOADC R, NUM | 1         |
| $P(M, T(N, T(N, R))) \Rightarrow$                            | $R \rightarrow M$                 | LOADM R, M   | 3         |
| $P(M, T(N, T(N, M)))$                                        |                                   |              | Total: 17 |

Obsérvese que, si asumimos por ahora que hay suficientes registros, la secuencia de instrucciones resultante en la tercera columna de la tabla si se lee en orden inverso (esto es, si se sigue el orden de instrucciones asociadas a las reglas de producción en orden de anti-derivación) y se hace una asignación correcta de registros nos da una traducción correcta de la expresión  $x+3*(7*y)$ :

```
LOADM R, M # y
LOADC R, NUM # 7
TIMESR R, R # 7*y
LOADC R, NUM # 3
TIMESR R, R # 3*(7*y)
LOADM R, M # x
PLUSR R, R # x+3*(7*y)
```

La gramática anterior es ambigua. El árbol de  $x+3*(7*y)$  puede ser generado también mediante la siguiente derivación a izquierdas:

| Otra derivación árbol a izquierdas para $P(M, T(N, T(N, M)))$ |                                                             |                  |           |
|---------------------------------------------------------------|-------------------------------------------------------------|------------------|-----------|
| Derivación                                                    | Producción                                                  | Instrucción      | Coste     |
| $R \Rightarrow$                                               | $R \rightarrow \text{PLUS}(R, \text{TIMES}(\text{NUM}, R))$ | PLUSCR R, NUM, R | 4         |
| $P(R, T(N, R)) \Rightarrow$                                   | $R \rightarrow M$                                           | LOADM R, M       | 3         |
| $P(M, T(N, R)) \Rightarrow$                                   | $R \rightarrow \text{TIMES}(R, M)$                          | TIMESM R, M      | 6         |
| $P(M, T(N, T(R, M)))$                                         | $R \rightarrow \text{NUM}$                                  | LOADC R, NUM     | 1         |
| $P(M, T(N, T(N, M)))$                                         |                                                             |                  | Total: 14 |

La nueva secuencia de instrucciones para  $x+3*(7*y)$  es:

```
LOADC R, NUM # 7
TIMESM R, M # 7*y
LOADM R, M # x
PLUSCR R, NUM, R # x+3*(7*y)
```

Cada antiderivación a izquierdas produce una secuencia de instrucciones que es una traducción legal del AST de  $x+3*(7*y)$ .

*El problema de la selección de código óptima puede aproximarse resolviendo el problema de encontrar la derivación árbol óptima que produce el árbol de entrada (en representación intermedia IR)*

**Definición 10.2.1.** *Un generador de generadores de código es una componente software que toma como entrada una especificación de la plataforma objeto -por ejemplo mediante una gramática árbol- y genera un módulo que es utilizado por el compilador. Este módulo lee la representación intermedia (habitualmente un árbol) y retorna código máquina como resultado.*

Un ejemplo de generador de generadores de código es iburg [?].

Véase también el libro Automatic Code Generation Using Dynamic Programming Techniques y la página <http://www.bytelabs.org/hburg.html>

**Ejercicio 10.2.1.** *Responda a las siguientes preguntas:*

- Sea  $G_M$  la gramática árbol asociada segun la descripción anterior con el juego de instrucciones de la máquina  $M$ . Especifique formalmente las cuatro componentes de la gramática  $G_M = (\Sigma_M, V_M, P_M, S_M)$
- ¿Cual es el lenguaje árbol generado por  $G_M$ ?
- ¿A que lenguaje debe pertenecer la representación intermedia IR para que se pueda aplicar la aproximación presentada en esta sección?

### 10.3. Patrones Árbol y Transformaciones Árbol

Una transformación de un programa puede ser descrita como un conjunto de *reglas de transformación* o *esquema de traducción árbol* sobre el árbol abstracto que representa el programa.

En su forma mas sencilla, estas reglas de transformación vienen definidas por ternas  $(p, e, action)$ , donde la primera componente de la terna  $p$  es un *patrón árbol* que dice que árboles deben ser seleccionados. La segunda componente  $e$  dice cómo debe transformarse el árbol que casa con el patrón  $p$ . La acción  $action$  indica como deben computarse los atributos del árbol transformado a partir de los atributos del árbol que casa con el patrón  $p$ . Una forma de representar este esquema sería:

$$p \Longrightarrow e \{ \text{action} \}$$

Por ejemplo:

$$PLUS(NUM_1, NUM_2) \Longrightarrow NUM_3 \{ \$NUM\_3\{VAL\} = \$NUM\_1\{VAL\} + \$NUM\_2\{VAL\} \}$$

cuyo significado es que dondequiera que haya un nódo del AAA que case con el *patrón de entrada*  $PLUS(NUM, NUM)$  deberá sustituirse el subárbol  $PLUS(NUM, NUM)$  por el subárbol  $NUM$ . Al igual que en los esquemas de traducción, enumeramos las apariciones de los símbolos, para distinguirlos en la parte semántica. La acción indica como deben recomputarse los atributos para el nuevo árbol: El atributo **VAL** del árbol resultante es la suma de los atributos **VAL** de los operandos en el árbol que ha casado. La transformación se repite hasta que se produce la *normalización del árbol*.

Las reglas de “casamiento” de árboles pueden ser mas complejas, haciendo alusión a propiedades de los atributos, por ejemplo

$$ASSIGN(LEFTVALUE, x) \text{ and } \{ \text{notlive}(\$LEFTVALUE\{VAL\}) \} \Longrightarrow NIL$$

indica que se pueden eliminar aquellos árboles de tipo asignación en los cuáles la variable asociada con el nodo  $LEFTVALUE$  no se usa posteriormente.

Otros ejemplos con variables  $S_1$  y  $S_2$ :

$$\begin{aligned} IFELSE(NUM, S_1, S_2) \text{ and } \{ \$NUM\{VAL\} \neq 0 \} &\Longrightarrow S_1 \\ IFELSE(NUM, S_1, S_2) \text{ and } \{ \$NUM\{VAL\} == 0 \} &\Longrightarrow S_2 \end{aligned}$$

Observe que en el patrón de entrada  $ASSIGN(LEFTVALUE, x)$  aparece un “comodín”: la variable-árbol  $x$ , que hace que el árbol patrón  $ASSIGN(LEFTVALUE, x)$  case con cualquier árbol de asignación, independientemente de la forma que tenga su subárbol derecho.

Las siguientes definiciones formalizan una aproximación simplificada al significado de los conceptos *patrones árbol* y *casamiento de árboles*.

#### Patrón Árbol

**Definición 10.3.1.** Sea  $(\Sigma, \rho)$  un alfabeto con función de aridad y un conjunto (puede ser infinito) de variables  $V = \{x_1, x_2, \dots\}$ . Las variables tienen aridad cero:

$$\rho(x) = 0 \quad \forall x \in V.$$

Un elemento de  $B(V \cup \Sigma)$  se denomina patrón sobre  $\Sigma$ .

#### Patrón Lineal

**Definición 10.3.2.** Se dice que un patrón es un patrón lineal si ninguna variable se repite.

**Definición 10.3.3.** Se dice que un patrón es de tipo  $(x_1, \dots, x_k)$  si las variables que aparecen en el patrón leídas de izquierda a derecha en el árbol son  $x_1, \dots, x_k$ .

**Ejemplo 10.3.1.** Sea  $\Sigma = \{A, CONS, NIL\}$  con  $\rho(A) = \rho(NIL) = 0, \rho(CONS) = 2$  y sea  $V = \{x\}$ . Los siguientes árboles son ejemplos de patrones sobre  $\Sigma$ :

$$\{x, \text{CONS}(A, x), \text{CONS}(A, \text{CONS}(x, \text{NIL})), \dots\}$$

El patrón  $\text{CONS}(x, \text{CONS}(x, \text{NIL}))$  es un ejemplo de patrón no lineal. La idea es que un patrón lineal como éste “fuerza” a que los árboles  $t$  que casen con el patrón deben tener iguales los dos correspondientes subárboles  $t/1$  y  $t/2$ .<sup>1</sup> situados en las posiciones de las variables

**Ejercicio 10.3.1.** Dado la gramática árbol:

$$\begin{aligned} S &\rightarrow S_1(a, S, b) \\ S &\rightarrow S_2(\text{NIL}) \end{aligned}$$

la cuál genera los árboles concretos para la gramática

$$S \rightarrow aSb \mid \epsilon$$

¿Es  $S_1(a, X(\text{NIL}), b)$  un patrón árbol sobre el conjunto de variables  $\{X, Y\}$ ? ¿Lo es  $S_1(X, Y, a)$ ? ¿Es  $S_1(X, Y, Y)$  un patrón árbol?

**Ejemplo 10.3.2.** Ejemplos de patrones para el AAA definido en el ejemplo 4.9.2 para el lenguaje Tutu son:

$$x, y, \text{PLUS}(x, y), \text{ASSIGN}(x, \text{TIMES}(y, \text{ID})), \text{PRINT}(y) \dots$$

considerando el conjunto de variables  $V = \{x, y\}$ . El patrón  $\text{ASSIGN}(x, \text{TIMES}(y, \text{ID}))$  es del tipo  $(x, y)$ .

## Sustitución

**Definición 10.3.4.** Una sustitución árbol es una aplicación  $\theta$  que asigna variables a patrones  $\theta : V \rightarrow B(V \cup \Sigma)$ .

Tal función puede ser naturalmente extendida de las variables a los árboles: los nodos (hoja) etiquetados con dichas variables son sustituidos por los correspondientes subárboles.

$$\begin{aligned} \theta &: B(V \cup \Sigma) \rightarrow B(V \cup \Sigma) \\ t\theta &= \begin{cases} x\theta & \text{si } t = x \in V \\ a(t_1\theta, \dots, t_k\theta) & \text{si } t = a(t_1, \dots, t_k) \end{cases} \end{aligned}$$

Obsérvese que, al revés de lo que es costumbre, la aplicación de la sustitución  $\theta$  al patrón se escribe por detrás:  $t\theta$ .

También se escribe  $t\theta = t\{x_1/x_1\theta, \dots, x_k/x_k\theta\}$  si las variables que aparecen en  $t$  de izquierda a derecha son  $x_1, \dots, x_k$ .

**Ejemplo 10.3.3.** Si aplicamos la sustitución  $\theta = \{x/A, y/\text{CONS}(A, \text{NIL})\}$  al patrón  $\text{CONS}(x, y)$  obtenemos el árbol  $\text{CONS}(A, \text{CONS}(A, \text{NIL}))$ . En efecto:

$$\text{CONS}(x, y)\theta = \text{CONS}(x\theta, y\theta) = \text{CONS}(A, \text{CONS}(A, \text{NIL}))$$

**Ejemplo 10.3.4.** Si aplicamos la sustitución  $\theta = \{x/\text{PLUS}(\text{NUM}, x), y/\text{TIMES}(\text{ID}, \text{NUM})\}$  al patrón  $\text{PLUS}(x, y)$  obtenemos el árbol  $\text{PLUS}(\text{PLUS}(\text{NUM}, x), \text{TIMES}(\text{ID}, \text{NUM}))$ :

$$\text{PLUS}(x, y)\theta = \text{PLUS}(x\theta, y\theta) = \text{PLUS}(\text{PLUS}(\text{NUM}, x), \text{TIMES}(\text{ID}, \text{NUM}))$$

---

<sup>1</sup>Repase la notación de Dewey introducida en la definición 4.9.6

## Casamiento Árbol

**Definición 10.3.5.** Se dice que un patrón  $\tau \in B(V \cup \Sigma)$  con variables  $x_1, \dots, x_k$  casa con un árbol  $t \in B(\Sigma)$  si existe una sustitución de  $\tau$  que produce  $t$ , esto es, si existen  $t_1, \dots, t_k \in B(\Sigma)$  tales que  $t = \tau\{x_1/t_1, \dots, x_k/t_k\}$ . También se dice que  $\tau$  casa con la sustitución  $\{x_1/t_1, \dots, x_k/t_k\}$ .

**Ejemplo 10.3.5.** El patrón  $\tau = \text{CONS}(x, \text{NIL})$  casa con el árbol  $t = \text{CONS}(\text{CONS}(A, \text{NIL}), \text{NIL})$  y con el subárbol  $t.1$ . Las respectivas sustituciones son  $t\{x/\text{CONS}(A, \text{NIL})\}$  y  $t.1\{x/A\}$ .

$$\begin{aligned} t &= \tau\{x/\text{CONS}(A, \text{NIL})\} \\ t.1 &= \tau\{x/A\} \end{aligned}$$

**Ejercicio 10.3.2.** Sea  $\tau = \text{PLUS}(x, y)$  y  $t = \text{TIMES}(\text{PLUS}(\text{NUM}, \text{NUM}), \text{TIMES}(\text{ID}, \text{ID}))$ . Calcule los subárboles  $t'$  de  $t$  y las sustituciones  $\{x/t_1, y/t_2\}$  que hacen que  $\tau$  case con  $t'$ .

Por ejemplo es obvio que para el árbol raíz  $t/\epsilon$  no existe sustitución posible:

$$t = \text{TIMES}(\text{PLUS}(\text{NUM}, \text{NUM}), \text{TIMES}(\text{ID}, \text{ID})) = \tau\{x/t_1, y/t_2\} = \text{PLUS}(x, y)\{x/t_1, y/t_2\}$$

ya que un término con raíz  $\text{TIMES}$  nunca podrá ser igual a un término con raíz  $\text{PLUS}$ .

El problema aquí es equivalente al de las expresiones regulares en el caso de los lenguajes lineales. En aquellos, los autómatas finitos nos proveen con un mecanismo para reconocer si una determinada cadena “casa” o no con la expresión regular. Existe un concepto análogo, el de *autómata árbol* que resuelve el problema del “casamiento” de patrones árbol. Al igual que el concepto de autómata permite la construcción de software para la búsqueda de cadenas y su posterior modificación, el concepto de autómata árbol permite la construcción de software para la búsqueda de los subárboles que casan con un patrón árbol dado.

## 10.4. El método m

El método `$yatl_object->m($t)` permite la búsqueda de los nodos que casan con una expresión regular árbol dada. En un contexto de lista devuelve una lista con nodos del tipo `Parse::Eyapp::Node::Match` que referencian a los nodos que han casado. Los nodos en la lista se estructuran según un árbol (atributos `children` y `father`) de manera que el padre de un nodo `$n` del tipo `Parse::Eyapp::Node::Match` es el nodo `$f` que referencia al inmediato antecesor en el árbol que ha casado.

Los nodos `Parse::Eyapp::Node::Match` son a su vez nodos (heredan de) `Parse::Eyapp::Node` y se estructuran según un árbol.

Los nodos aparecen en la lista retornada en orden primero profundo de recorrido del árbol `$t`.

Un nodo `$r` de la clase `Parse::Eyapp::Node::Match` dispone de varios atributos y métodos:

- El método `$r->node` retorna el nudo del árbol `$t` que ha casado
- El método `$r->father` retorna el nodo padre en el árbol the matching. Se cumple que `$r->father->node` es una referencia al antepasado mas cercano en `$t` de `$r->node` que casa con el patrón árbol (`$SimpleTrans::blocks` en el ejemplo que sigue).
- El método `$r->coord` retorna las coordenadas del nudo del árbol que ha casado usando una notación con puntos. Por ejemplo la coordenada `".1.3.2"` denota al nodo `$t->child(1)->child(3)->child(2)`, siendo `$t` la raíz del árbol de búsqueda.
- El método `$r->depth` retorna la profundidad del nudo del árbol que ha casado

La clase `Parse::Eyapp::Node::Match` hereda de `Parse::Eyapp::Node`. Dispone además de otros métodos para el caso en que se usan varios patrones en la misma búsqueda.

Por ejemplo, método `$r->names` retorna una referencia a los nombres de los patrones que han casado con el nodo.

En un contexto escalar `m` devuelve el primer elemento de la lista de nodos `Parse::Eyapp::Node::Match`.

## El método m: Ejemplo de uso

La siguiente expresión regular árbol casa con los nodos *bloque*:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> sed -ne '2p' Trans.trg
blocks: /BLOCK|FUNCTION|PROGRAM/
```

Utilizaremos el siguiente programa de entrada:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/script> cat -n blocks.c
1 test (int n)
2 {
3 while (1) {
4 if (1>0) {
5 int a;
6 break;
7 }
8 else if (2> 0){
9 char b;
10 continue;
11 }
12 }
13 }
```

Ejecutamos la versión con análisis de tipos de nuestro compilador de simple C. SimpleC es una versión simplificada de C (véase el capítulo 12 y las secciones 12.3 y 12.4), bajo el control del depurador:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/script> perl -wd usetypes.pl blocks.c
Loading DB routines from perl5db.pl version 1.28
Editor support available.
Enter h or 'h h' for help, or 'man perldebug' for more help.
main:.(usetypes.pl:6): my $filename = shift || die "Usage:\n$0 file.c\n";
DB<1> b Simple::Types::compile
DB<2> c
Simple::Types::compile(Types.eypp:529):
529: my($self)=shift;
DB<2> l 540,545
540: our $blocks;
541: my @blocks = $blocks->m($t);
542: $_->node->{fatherblock} = $_->father->{node} for (@blocks[1..$#blocks]);
543
544 # Type checking
545: set_types($t);
DB<3> c 545
Simple::Types::compile(Types.eypp:545):
545: set_types($t);
```

Nos detenemos en la línea 545 justo después de la llamada a `m` (línea 541) y de actualizar el atributo `fatherblock` de los nodos bloque. La línea 542 crea una jerarquía árbol con los nodos bloque que se superpone a la jerarquía del árbol sintáctico abstracto. En esta línea `$_->father->{node}` es el nodo del AST que es un ancestro de `$_->node` y que ha casado con la expresión regular árbol.

Este sub-árbol de bloques puede ayudarnos en la fase de ubicar que declaración se aplica a una aparición de un objeto dado en el texto del programa. Para cada ocurrencia debemos determinar si fué declarada en el bloque en el que ocurre o bien en uno de los nodos bloque antepasados de este.

```
DB<4> $Parse::Eyapp::Node::INDENT = 2
DB<5> x $blocks[0]->str
```

```

0 '
Match[[PROGRAM:0:blocks]](
 Match[[FUNCTION:1:blocks:test]](
 Match[[BLOCK:6:blocks:8:4:test]],
 Match[[BLOCK:5:blocks:4:4:test]]
)
) # Parse::Eyapp::Node::Match'

```

La orden `x $blocks[0]->str` nos permite visualizar el árbol de casamientos. El método `coord` nos devuelve una cadena con las coordenadas del nodo que ha casado. Las coordenadas son una secuencia de puntos y números que describe el camino para llegar al nodo desde la raíz del árbol. Veamos las coordenadas de los nodos en el AST del ejemplo:

```

DB<6> x map {$_->coord} @blocks
0 ''
1 '.0'
2 '.0.0.1.0.1'
3 '.0.0.1.0.2.1'
DB<7> x $t->child(0)->child(0)->child(1)->child(0)->child(2)->child(1)->str
0 '
BLOCK[8:4:test]~{0}(
 CONTINUE[10,10]
)
DB<8> x $t->descendant('.0.0.1.0.2.1')->str
0 '
BLOCK[8:4:test]~{0}(
 CONTINUE[10,10]
)
DB<9> x $t->descendant('.0.0.1.0.1')->str
0 '
BLOCK[4:4:test]~{0}(
 BREAK[6,6]
)

```

El método `descendant` es un método de los objetos `Parse::Eyapp::Node` y retorna una referencia al nodo descrito por la cadena de coordenadas que se le pasa como argumento. En este sentido el método `child` es una especialización de `descendant`.

La profundidad de los nodos que han casado puede conocerse con el método `depth`:

```

DB<9> x map {$_->depth} @blocks
0 0
1 1
2 5
3 6
DB<16> x map {ref($_->node)} @blocks
0 'PROGRAM'
1 'FUNCTION'
2 'BLOCK'
3 'BLOCK'
DB<17> x map {ref($_->father)} @blocks
0 ''
1 'Parse::Eyapp::Node::Match'
2 'Parse::Eyapp::Node::Match'
3 'Parse::Eyapp::Node::Match'
DB<19> x map {ref($_->father->node)} @blocks[1..3]

```



```

0 'PROGRAM'
1 'FUNCTION'
2 'FUNCTION'

```

## 10.5. Transformaciones Arbol con treereg

La distribución de Eyapp provee un conjunto de módulos que permiten la manipulación y transformación del AST. La forma más usual de hacerlo es escribir el programa de transformaciones árbol en un fichero separado. El tipo `.trg` se asume por defecto. El ejemplo que sigue trabaja en cooperación con el programa Eyapp presentado en la sección 8.11 página 8.11. El programa árbol sustituye nodos producto por un número potencia de dos por nodos unarios de desplazamiento a izquierda. De este modo facilitamos la tarea de la generación de un código más eficiente:

```

pl@nereida:~/LEyapp/examples$ cat -n Shift.trg
1 # File: Shift.trg
2 {
3 sub log2 {
4 my $n = shift;
5 return log($n)/log(2);
6 }
7
8 my $power;
9 }
10 mult2shift: TIMES($e, NUM($m))
11 and { $power = log2($m->{attr}); (1 << $power) == $m->{attr} } => {
12 $_[0]->delete(1);
13 $_[0]->{shift} = $power;
14 $_[0]->type('SHIFTLEFT');
15 }

```

Obsérvese que la variable léxica `$power` es visible en la definición de la transformación `mult2shift`. Compilamos el programa de transformaciones usando el guión `treereg` :

```

nereida:~/src/perl/YappWithDefaultAction/examples> treereg Shift
nereida:~/src/perl/YappWithDefaultAction/examples> ls -ltr | tail -1
-rw-rw---- 1 pl users 1405 2006-11-06 14:09 Shift.pm

```

Como se ve, la compilación produce un módulo que `Shift.pm` que contiene el código de los analizadores. El código generado por la versión 0.2 de `treereg` es el siguiente:

```

pl@nereida:~/LEyapp/examples$ cat -n Shift.pm
1 package Shift;
2
3 # This module has been generated using Parse::Eyapp::Treereg
4 # from file Shift.trg. Don't modify it.
5 # Change Shift.trg instead.
6 # Copyright (c) Casiano Rodriguez-Leon 2006. Universidad de La Laguna.
7 # You may use it and distribute it under the terms of either
8 # the GNU General Public License or the Artistic License,
9 # as specified in the Perl README file.
10
11 use strict;
12 use warnings;
13 use Carp;
14 use Parse::Eyapp::_TreeregexpSupport qw(until_first_match checknumchildren);

```

```

15
16 our @all = (our $mult2shift,) = Parse::Eyapp::YATW->buildpatterns(mult2shift => \&mult2s
17
18 #line 2 "Shift.trg"
19
20 sub log2 {
21 my $n = shift;
22 return log($n)/log(2);
23 }
24
25 my $power;
26
27
28 sub mult2shift {
29 my $mult2shift = $_[3]; # reference to the YATW pattern object
30 my $e;
31 my $NUM;
32 my $TIMES;
33 my $m;
34
35 {
36 my $child_index = 0;
37
38 return 0 unless (ref($TIMES = $_[$child_index]) eq 'TIMES');
39 return 0 unless defined($e = $TIMES->child(0+$child_index));
40 return 0 unless (ref($NUM = $TIMES->child(1+$child_index)) eq 'NUM');
41 return 0 unless defined($m = $NUM->child(0+$child_index));
42 return 0 unless do
43 #line 10 "Shift.trg"
44 { $power = log2($m->{attr}); (1 << $power) == $m->{attr} };
45
46 } # end block of child_index
47 #line 11 "Shift.trg"
48 {
49 $_[0]->delete(1);
50 $_[0]->{shift} = $power;
51 $_[0]->type('SHIFTLEFT');
52 }
53
54 } # end of mult2shift
55
56
57 1;

```

## Uso de los Módulos Generados

El programa de análisis se limita a cargar los dos módulos generados, crear el analizador, llamar al método Run para crear el árbol y proceder a las sustituciones mediante la llamada `$t->s(@Shift::all):`

```

pl@nereida:~/LEyapp/examples$ cat -n useruleandshift.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Rule5;
4 use Parse::Eyapp::Base qw(insert_function);
5 use Shift;

```

```

6
7 sub SHIFTLEFT::info { $_[0]{shift} }
8 insert_function('TERMINAL::info', \&TERMINAL::attr);
9
10 my $parser = new Rule5();
11 my $t = $parser->Run;
12 print "*****\n",$t->str,"\n";
13 $t->s(@Shift::all);
14 print "*****\n",$t->str,"\n";

```

Cuando alimentamos el programa con la entrada `a=b*32` obtenemos la siguiente salida:

```

pl@nereida:~/LEyapp/examples$ useruleandshift.pl
a=b*32

ASSIGN(TERMINAL[a],TIMES(VAR(TERMINAL[b]),NUM(TERMINAL[32])))

ASSIGN(TERMINAL[a],SHIFTLEFT[5](VAR(TERMINAL[b])))

```

## 10.6. Transformaciones de Árboles con `Parse::Eyapp::Treeregexp`

El módulo `Parse::Eyapp::Treeregexp` permite la transformación de árboles mediante el uso de *Expresiones Regulares Arbol*. Las expresiones regulares árbol serán introducidas en mas detalle en la sección 4.9.

### Optimización del Traductor de Infijo a Postfijo

El siguiente ejemplo modifica el anterior esquema de traducción de infijo a postfijo para producir un código de postfijo mas eficiente. Para ello se transforma el árbol generado durante la fase *Tree Construction Time* y antes de la fase *Execution Time*. El código *Treeregexp* que define el conjunto de transformaciones se encuentra en las líneas 74-103.

Las transformaciones consisten en

1. Compactar los árboles `UMINUS` a un número negativo
2. Realizar *plegado de constantes*: sustituir los árboles de constantes por su evaluación
3. Sustituir los árboles producto en los que uno de los factores es cero por el número cero.

Después de ello se realiza la traducción quedando la misma como el atributo `t` del nodo raíz (línea 120).

A partir de este momento, si el traductor tuviera un mayor número de fases de posterior tratamiento del árbol, los nodos de tipo código y los nodos hoja cuya funcionalidad es puramente sintáctica como los terminales `=`, `*` etc. pueden ser eliminados. Es por eso que los suprimimos en las líneas 122-123.

Veamos primero el código y luego lo discutiremos en mas detalle:

```

nereida:~/src/perl/YappWithDefaultAction/examples> cat -n TSwithtreetransformations.eyp
1 # File TSwithtreetransformations.eyp
2 %right '='
3 %left '- ' '+'
4 %left '* ' '/'
5 %left NEG
6
7 %{
8 # Treeregexp is the engine for tree transformations
9 use Parse::Eyapp::Treeregexp;

```

```

10 use Data::Dumper;
11 $Data::Dumper::Indent = 1;
12 $Data::Dumper::Deepcopy = 1;
13 $Data::Dumper::Deparse = 1;
14 %}
15
16 %metatree
17
18 %defaultaction {
19 if (@_==4) { # binary operations: 4 = lhs, left, operand, right
20 $lhs->{t} = "$_[1]->{t} $_[3]->{t} $_[2]->{attr}";
21 return
22 }
23 die "Fatal Error. Unexpected input\n".Dumper(@_);
24 }
25
26 %%
27 line: %name PROG
28 exp <%name EXP + ','>
29 { @{$lhs->{t}} = map { $_->{t}} ($lhs->child(0)->Children()); }
30
31 ;
32
33 exp: %name NUM NUM { $lhs->{t} = $_[1]->{attr}; }
34 | %name VAR VAR { $lhs->{t} = $_[1]->{attr}; }
35 | %name ASSIGN VAR '=' exp { $lhs->{t} = "$_[1]->{attr} $_[3]->{t} =" }
36 | %name PLUS exp '+' exp
37 | %name MINUS exp '-' exp
38 | %name TIMES exp '*' exp
39 | %name DIV exp '/' exp
40 | %name UMINUS '-' exp %prec NEG { $_[0]->{t} = "$_[2]->{t} NEG" }
41 | '(' exp ')' %begin { $_[2] } /* skip parenthesis */
42 ;
43
44 %%
45
46 # subroutines _Error and _Lexer
47
48
49 sub Run {
50 my($self)=shift;
51 print "input: "; $x = <>;
52 my $tree = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
53 #yydebug => 0xFF
54);
55
56 my $transform = Parse::Eyapp::Treeregexp->new(STRING => q{
57
58 delete_code : CODE => { $delete_code->delete() }
59
60 {
61 sub not_semantic {
62 my $self = shift;

```

```

81 return 1 if $self->{token} eq $self->{attr};
82 return 0;
83 }
84 }
85
86 delete_tokens : TERMINAL and { not_semantic($TERMINAL) } => { $delete_tokens->dele
87
88 delete = delete_code delete_tokens;
89
90 uminus: UMINUS(., NUM($x), .) => { $x->{attr} = -$x->{attr}; $_[0] = $NUM }
91
92 constantfold: /TIMES|PLUS|DIV|MINUS/(NUM($x), ., NUM($y))
93 => {
94 $x->{attr} = eval "$x->{attr} $W->{attr} $y->{attr}";
95 $_[0] = $NUM[0];
96 }
97
98 zero_times: TIMES(NUM($x), ., .) and { $x->{attr} == 0 } => { $_[0] = $NUM }
99 times_zero: TIMES(., ., NUM($x)) and { $x->{attr} == 0 } => { $_[0] = $NUM }
100
101 algebraic_transformations = constantfold zero_times times_zero;
102
103 },
104 PACKAGE => 'TsWithtreetransformations',
105 OUTPUTFILE => 'main.pm',
106 SEVERITY => 0,
107 NUMBERS => 0,
108);
109
110 # Create the transformer
111 $transform->generate();
112 # Get the AST
113
114 our ($uminus);
115 $uminus->s($tree);
116
117 our (@algebraic_transformations);
118 $tree->s(@algebraic_transformations);
119
120 $tree->translation_scheme();
121
122 our (@delete);
123 $tree->s(@delete);
124 print Dumper($tree);
125 }

```

### La Estructura de un Programa Treeregexp

La estructura de un programa Treeregexp es sencilla. Consiste en la repetición de tres tipos de expresiones regulares árbol: las *treeregexp* propiamente dichas, *código auxiliar para las transformaciones* y *definiciones de familias de transformaciones*.

```

treeregexplist:
 treeregexp*
;

```

```

treeregexp:
 IDENT ':' treereg ('=>' CODE)? # Treeregexp
| CODE # Código auxiliar
| IDENT '=' IDENT + ';' # Familia de transformaciones
;

```

Las expresiones regulares árbol propiamente dichas siguen la regla

```
IDENT ':' treereg ('=>' CODE)?
```

Podemos ver ejemplos de instancias de esta regla en las líneas 76, 86, 90, 92-96, 98 y 99. El identificador IDENT da el nombre a la regla. Actualmente (2006) existen estos tipos de treereg :

```

treereg:
 /* patrones con hijos */
 IDENT '(' childlist ')' ('and' CODE)?
| REGEXP ':' IDENT)? '(' childlist ')' ('and' CODE)?
| SCALAR '(' childlist ')' ('and' CODE)?
| '.' '(' childlist ')' ('and' CODE)?
 /* hojas */
| IDENT ('and' CODE)?
| REGEXP ':' IDENT)? ('and' CODE)?
| '.' ('and' CODE)?
| SCALAR ('and' CODE)?
| ARRAY
| '*'

```

## Las Reglas Treeregexp

Una regla como

```
zero_times: TIMES(NUM($x), ., .) and { $x->{attr} == 0 } => { $_[0] = $NUM }
```

crea un objeto transformación (concretamente un objeto de la clase Parse::Eyapp:YATW ) que puede ser referido a través de la variable escalar \$zero\_times. La primera parte de la regla zero\_times indica que para que se produzca el emparejamiento es necesario que el nodo visitado sea del tipo TIMES y su primer hijo es de tipo NUM. Una referencia al nodo hijo de NUM será automáticamente guardada en la variable léxica \$x.

## Escalares

El efecto de un escalar en una treeregexp es casar con cualquier nodo y almacenar su referencia en la variable.

La aparición de \$x en la treeregexp anterior casará con cualquier nodo. La referencia al nodo que ha casado queda en \$x. Así \$x podrá ser usado en el *patrón árbol semántico* o *condición semántica* (esto es, en el código opcional que va precedido de la palabra reservada **and**) y en la *acción de transformación árbol* (el código opcional que va precedido de la flecha gorda =>).

## El Punto

Un punto también casa con cualquier nodo. Puede verse como una abreviación de la expresión regular árbol escalar. Una referencia al nodo que casa queda almacenada en la variable léxica especial \$W . Si la expresión regular árbol tiene varios puntos sus referencias quedan almacenadas en la variable array @W. Es un error usar el identificador W en una expresión regular escalar. Por ejemplo, una treeregexp como:

```
constantfold: /TIMES|PLUS|DIV|MINUS/(NUM($W), ., NUM($y))
```

da lugar al error:

```
Error Can't use $W to identify an scalar treeregexp at line 100.
```

### Condiciones Semánticas

La segunda parte de la regla es opcional y comienza con la palabra reservada **and** seguida de un código que *explicita las condiciones semánticas que debe cumplir el nodo para que se produzca el casamiento*. En el ejemplo se explicita que el atributo del nodo (forzosamente del tipo **TERMINAL** en este caso) referenciado por **\$x** debe ser cero.

### Referenciado de los Nodos del Arbol

Es posible dentro de las partes de código referirse a los nodos del árbol. Cuando la rutina de transformación generada por el compilador para una **treeregexp** es llamada, el primer argumento **\$\_[0]** contiene la referencia al nodo que esta siendo visitado.

**Parse::Eyapp::Treeregexp** crea variables léxicas con nombres los tipos de los nodos a los que referencian. Así la subrutina generada para la transformación **zero\_times**

```
zero_times: TIMES(NUM($x), ., .) and { $x->{attr} == 0 } => { $_[0] = $NUM }
```

guarda en la variable lexic **\$TIMES** una copia de **\$\_[0]** y en la variable léxica **\$NUM** una referencia al nodo **\$TIMES->child(0)**.

Si un tipo de nodo se repite en la **treeregexp** la variable léxica asociada con dicho tipo se autodeclara como un array. Este es el caso de la transformación **constantfold** en la cual aparecen dos nodos de tipo **NUM**:

```
92 constantfold: /TIMES|PLUS|DIV|MINUS/(NUM($x), ., NUM($y))
93 => {
94 $x->{attr} = eval "$x->{attr} $W->{attr} $y->{attr}";
95 $_[0] = $NUM[0];
96 }
```

La variable **@NUM** es automáticamente declarada: **\$NUM[0]** es una referencia al primer nodo **NUM** y **\$NUM[1]** es una referencia al segundo.

### Código de Transformación

La tercera parte de la regla es también opcional y viene precedida de la *flecha gorda*. Habitualmente contiene el código que transforma el árbol. Para lograr la modificación del nodo del árbol visitado el programador **Treeregexp** deberá usar **\$\_[0]**. Recuerde que los elementos en **@\_** son alias de los argumentos. Si el código de la tercera parte fuera reescrito como:

```
{ $TIMES = $NUM }
```

no funcionaría ya que estaríamos modificando la variable léxica que referencia al nodo raíz del subarbol que ha casado.

### Expresiones Regulares

Es posible usar una *expresión regular lineal* alias *expresión regular clásica* alias *regexp* para explicitar el tipo de un nodo como indica la regla de producción:

```
treereg: REGEXP (':' IDENT)? '(' childlist ')' ('and' CODE)?
```

La **treeregexp** para el plegado de constantes constituye un ejemplo:

```

92 constantfold: /TIMES|PLUS|DIV|MINUS/(NUM($x), ., NUM($y))
93 => {
94 $x->{attr} = eval "$x->{attr} $W->{attr} $y->{attr}";
95 $_[0] = $NUM[0];
96 }

```

La expresión regular deberá especificarse entre barras de división (/) y es posible especificar opciones después del segundo slash (e, i, etc.). El identificador opcional después de la regexp indica el nombre para la variable léxica que almacenará una copia de la referencia al nodo del árbol. En el ejemplo `$bin` podría usarse para referenciar al nodo apuntado por `$_[0]`. Si no se especifica identificador quedará almacenado en la variable léxica especial `$W`. Si la expresión regular árbol tiene varias regexp (y/o puntos) sus referencias quedan almacenadas en la variable array `@W`.

La semántica de las *expresiones regulares lineales* es modificada ligeramente por `Parse::Eyapp::Treeregexp`. Por defecto se asume la opción `x`. El compilador de expresiones regulares árbol procede a la *inserción automática de la opción x*. Use la nueva opción `X` (`X` mayúscula) para evitar esta conducta. *Tampoco es necesario añadir anclas de frontera de palabra* `\b` a los identificadores que aparecen en la expresión regular lineal: el compilador de expresiones regulares árbol las insertará. Use la nueva opción `B` (`B` mayúscula) para negar esta conducta.

El siguiente ejemplo tomado del análisis de tipos en el compilador de Simple C muestra como no es necesario usar `x`:

```

67 # Binary Operations
68 bin: / PLUS
69 | MINUS
70 | TIMES
71 | DIV
72 | MOD
73 | GT
74 | GE
75 | LE
76 | EQ
77 | NE
78 | LT
79 | AND
80 | EXP
81 | OR
82 /($x, $y)
83 => {
84 $x = char2int($_[0], 0);
85 $y = char2int($_[0], 1);
86
87 if (($x->{t} == $INT) and ($y->{t} == $INT)) {
88 $_[0]->{t} = $INT;
89 return 1;
90 }
91 type_error("Incompatible types with operator '".($_[0]->lexeme)."', $_[0]->line);
92 }

```

Con la semántica habitual de las regexp la palabra reservada `WHILE` casaría con la subexpresión `LE` en la línea 76 provocando un análisis de tipos erróneo para esta clase de nodos. La *inserción automática de anclas* por parte de `Parse::Eyapp::Treeregexp` evita este - normalmente indeseable - efecto.

## Familias de Transformaciones

Las transformaciones creadas por `Parse::Eyapp::Treeregexp` pueden agruparse por familias. Esta es la función de la regla de producción:



```
treeregexp: IDENT '=' IDENT + ';' ;
```

En el ejemplo creamos una nueva familia denominada `algebraic_transformations` mediante la asignación de la línea 101:

```
algebraic_transformations = constantfold zero_times times_zero;
```

Las transformaciones en esa familia pueden ser accedidas posteriormente para su aplicación mediante la variable de paquete `@algebraic_transformations` (véanse las líneas 117-118).

## Código de Apoyo

En medio de la definición de cualquier regla `treeregexp` es posible insertar código de apoyo siempre que se sitúe entre llaves:

```
78 {
79 sub not_semantic {
80 my $self = shift;
81 return 1 if $self->{token} eq $self->{attr};
82 return 0;
83 }
84 }
85
86 delete_tokens : TERMINAL and { not_semantic($TERMINAL) } => { $delete_tokens->delete
```

## El método delete de los objetos YATW

Los objetos de la clase `Parse::Eyapp::YATW` como `$delete_tokens` disponen de un método `delete` que permite eliminar con seguridad un hijo de la raíz del subárbol que ha casado. En este caso los nodos que casan son los de la clase `TERMINAL` en los que el valor de la clave `token` coincide con el valor de la clave `attr`.

## Fases en la Ejecución de un Programa Treeregexp

Un programa `Treeregexp` puede - como en el ejemplo - proporcionarse como una cadena de entrada al método `new` de la clase `Parse::Eyapp::Treeregexp` o bien escribirse en un fichero separado (la extensión `.trg` es usada por defecto) y compilado con el script `treereg` que acompaña a la distribución de `Parse::Eyapp`.

La primera fase en la ejecución de un programa `Treeregexp` es la *fase de creación del paquete Treeregexp* que contendrá las subrutinas de reconocimiento de los patrones árbol definidos en el programa `Treeregexp`. En el ejemplo esta fase tiene lugar en las líneas 74-111 con las llamadas a `new` (que crea el objeto programa `Treeregexp`) y `generate` (que crea el paquete conteniendo las subrutinas reconocedoras).

```
74 my $transform = Parse::Eyapp::Treeregexp->new(STRING => q{
75
76 delete_code : CODE => { $delete_code->delete() }
...
103 },
104 PACKAGE => 'TsWithtreetransformations',
105 OUTPUTFILE => 'main.pm',
106 SEVERITY => 0,
107 NUMBERS => 0,
108);
109
110 # Create the transformer
111 $transform->generate();
```

Durante esta fase se crea un *objeto transformación* i(perteneciente a la clase `Parse::Eyapp::YATW`) por cada expresión regular árbol que aparece en el programa `Treeregexp`. Las variables contenedor de cada uno de esos objetos tienen por nombre el que se les dió a las correspondientes expresiones regulares árbol. En nuestro ejemplo, y después de esta fase habrán sido creadas variables escalares de paquete `$delete_tokens`, `$delete`, `$uminus`, `$constantfold`, `$zero_times` y `$times_zero` asociadas con cada una de las expresiones regulares árbol definidas. También se crearán variables array para cada una de las familias de transformaciones especificadas: Así la variable `@delete` contiene (`$delete_tokens`, `$delete`) y la variable `@algebraic_transformations` es igual a (`$constantfold`, `$zero_times`). La variable de paquete especial `@all` es un array que contiene todas las transformaciones definidas en el programa.

Una vez creados los objetos transformación y las familias de transformaciones `Parse::Eyapp::YATW` podemos proceder a transformar el árbol mediante su uso. Esto puede hacerse mediante el método `s` el cual procede a modificar el árbol pasado como parámetro.

```
114 our ($uminus);
115 $uminus->s($tree);
```

En el ejemplo la llamada `$uminus->s($tree)` da lugar al recorrido primero-profundo de `$tree`. Cada vez que un nodo casa con la regla:

```
UMINUS(., NUM($x), .) # first child is '-', second the number and third the code
```

se le aplica la transformación:

```
{ $x->{attr} = -$x->{attr}; $_[0] = $NUM }
```

Esta transformación hace que el nodo `UMINUS` visitado sea sustituido por un nodo de tipo `NUM` cuyo atributo sea el número de su hijo `NUM` cambiado de signo.

Las líneas 117-118 nos muestran como someter un árbol a un conjunto de transformaciones:

```
117 our (@algebraic_transformations);
118 $tree->s(@algebraic_transformations);
```

Los objetos nodo (esto es, los que pertenecen a la clase `Parse::Eyapp::Node`) disponen del método `s` que recibe como argumentos una familia de transformaciones. La familia de transformaciones es aplicada iterativamente al árbol hasta que este no cambia.

Nótese que consecuencia de esta definición es que es posible escribir transformaciones que dan lugar a bucles infinitos. Por ejemplo si en `@algebraic_transformations` incluimos una transformación que aplique las propiedades conmutativas de la suma:

```
commutative_add: PLUS($x, ., $y, .)
=> { my $t = $x; $_[0]->child(0, $y); $_[0]->child(2, $t) }
```

el programa podría caer en un bucle infinito ya que la transformación es susceptible de ser aplicada indefinidamente. Sin embargo no se produce bucle infinito si llamamos al código asociado a la transformación:

```
$commutative_add->($tree);
```

ya que en este caso se produce un sólo recursivo descendente aplicando la transformación `$commutative_add`.

El uso de transformaciones conmutativas no tiene porque dar lugar a la no finalización del programa. La parada del programa se puede garantizar si *podemos asegurar que la aplicación reiterada del patrón implica la desaparición del mismo*. Por ejemplo, la transformación `comasocfold` puede ser añadida a la familia `algebraic_transformations` sin introducir problemas de parada:

```
comasocfold: TIMES(DIV(NUM($x), ., $b), ., NUM($y))
=> {
 $x->{attr} = $x->{attr} * $y->{attr};
```

```

 $_[0] = $DIV;
}

```

algebraic\_transformations = constantfold zero\_times times\_zero comasocfold;

La introducción de esta transformación permite el plegado de entradas como  $a=2;b=2/a*3$ :

```

nereida:~/src/perl/YappWithDefaultAction/examples> usetswithtreetransformations3.pl
a=2;b=2/a*3
$VAR1 = bless({ 'children' => [
 bless({ 'children' => [
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL'),
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => 2, 'token' => 'NUM' }, 'TERMINAL')],
 't' => 2
 }, 'NUM')
],
 't' => 'a 2 ='
 }, 'ASSIGN'),
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => 'b', 'token' => 'VAR' }, 'TERMINAL'),
 bless({ 'children' => [
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => 6, 'token' => 'NUM' }, 'TERMINAL')
],
 't' => 6
 }, 'NUM'),
 bless({ 'children' => [
 bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL')
],
 't' => 'a'
 }, 'VAR')
],
 't' => '6 a /'
}, 'DIV')
],
 't' => 'b 6 a / ='
}, 'ASSIGN')
]
}, 'EXP')
],
't' => ['a 2 =', 'b 6 a / =']
}, 'PROG');

```

## Ejecución del Ejemplo

Una vez compilado el analizador, escribimos el programa que usa el módulo generado:

```

nereida:~/src/perl/YappWithDefaultAction/examples> cat -n usetswithtreetransformations.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use TSwithtreetransformations;
4 use Parse::Eyapp::Treeregexp;
5
6 my $parser = TSwithtreetransformations->new();

```

7 \$parser->Run;

Al ejecutarlo obtenemos la siguiente salida:

```
nereida:~/src/perl/YappWithDefaultAction/examples> eyapp TSwithtreetransformations.eyyp ; \\
 usetswithtreetransformations.pl
```

input: a=2\*-3;b=a\*(2-1-1);c=a+b

```
$VAR1 = bless({ 'children' => [
 bless({ 'children' => [
 | bless({ 'children' => [
 | | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL'),
 | | bless({ 'children' => [
 | | | bless({ 'children' => [], 'attr' => -6, 'token' => 'NUM' }, 'TERMINAL')
 | | |],
 | | | 't' => -6
 | | }, 'NUM')
 | |],
 | 't' => 'a -6 ='
 | }, 'ASSIGN'),
 | bless({ 'children' => [
 | | bless({ 'children' => [], 'attr' => 'b', 'token' => 'VAR' }, 'TERMINAL'),
 | | bless({ 'children' => [
 | | | bless({ 'children' => [], 'attr' => 0, 'token' => 'NUM' }, 'TERMINAL')
 | | |],
 | | | 't' => 0
 | | }, 'NUM')
 | |],
 | 't' => 'b 0 ='
 | }, 'ASSIGN'),
 | bless({ 'children' => [
 | | bless({ 'children' => [], 'attr' => 'c', 'token' => 'VAR' }, 'TERMINAL'),
 | | bless({ 'children' => [
 | | | bless({ 'children' => [
 | | | | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL')
 | | | |],
 | | | | 't' => 'a'
 | | | | }, 'VAR'),
 | | | | bless({ 'children' => [
 | | | | | bless({ 'children' => [], 'attr' => 'b', 'token' => 'VAR' }, 'TERMINAL')
 | | | | |],
 | | | | | 't' => 'b'
 | | | | | }, 'VAR')
 | | | |],
 | | | | 't' => 'a b +'
 | | | }, 'PLUS')
 | |],
 | | 't' => 'c a b + ='
 | }, 'ASSIGN')
 |]
 }, 'EXP')
],
't' => ['a -6 =', 'b 0 =', 'c a b + =']
}, 'PROG');
```

Como puede verse la traducción de la frase de entrada a=2\*-3;b=a\*(2-1-1);c=a+b queda como

atributo `t` del nodo raíz `PROG`.

### Expresiones Regulares Arbol Array

Una expresión regular árbol array se escribe insertando un array Perl en la expresión regular árbol, por ejemplo `@a`. Una *expresión regular árbol array* `@a` casa con la secuencia mas corta de hijos del nodo tal que la siguiente expresión regular árbol (no array) casa. La lista de nodos que han casado con la expresión regular árbol array quedará en la variable léxica `@a`. Por ejemplo, después de un casamiento de un árbol `$t` con la expresión regular árbol `BLOCK(@a, ASSIGN($x, $e), @b)`, la variable léxica `@a` contendrá la lista de nodos hijos de `$t` que precede a la primera aparición de `ASSIGN($x, $e)`. Si no existe expresión regular árbol siguiente - el caso de `@b` en el ejemplo - la expresión regular array casará con todos los nodos hijo a partir del último casamiento (no array). Así `@b` contendrá la lista de referencias a los nodos hijos de `$t` posteriores al nodo `ASSIGN($x, $e)`.

Es ilegal escribir dos expresiones regulares arbol array seguidas (por ejemplo `A(@a, @b)`).

El siguiente ejemplo muestra como usar las expresiones árbol array para mover las asignaciones invariantes de un bucle fuera del mismo (líneas 104-116):

```
nereida:~/src/perl/YappWithDefaultAction/t> cat -n 34moveinvariantoutofloopcomplexformula.t
 1 #!/usr/bin/perl -w
 2 use strict;
 5 use Parse::Eyapp;
 6 use Data::Dumper;
 7 use Parse::Eyapp::Treeregexp;
 8
 9 my $grammar = q{
10 %{
11 use Data::Dumper;
12 %}
13 %right '='
14 %left '-' '+'
15 %left '*' '/'
16 %left NEG
17 %tree
18
19 %%
20 block: exp <%name BLOCK + ';'> { $_[1] }
21 ;
22
23 exp: %name NUM
24 NUM
25 | %name WHILE
26 'while' exp '{' block '}'
27 | %name VAR
28 VAR
29 | %name ASSIGN
30 VAR '=' exp
31 | %name PLUS
32 exp '+' exp
33 | %name MINUS
34 exp '-' exp
35 | %name TIMES
36 exp '*' exp
37 | %name DIV
38 exp '/' exp
39 | %name UMINUS
```

```

40 '-' exp %prec NEG
41 | '(' exp ')' { $_[2] } /* Let us simplify a bit the tree */
42 ;
43
44 %%
45
87 }; # end grammar
88
46
99 $parser->YYData->{INPUT} = "a =1000; c = 1; while (a) { c = c*a; b = 5; a = a-1 }\n";
100 my $t = $parser->Run;
101 print "\n***** Before *****\n";
102 print Dumper($t);
104 my $p = Parse::Eyapp::Treeregexp->new(STRING => q{
105 moveinvariant:
106 BLOCK(@prests, WHILE(VAR($b), BLOCK(@a, ASSIGN($x, $e), @c)), @possts)
107 and { is_invariant($ASSIGN, $WHILE) } /* Check if ASSIGN is invariant relative */
108 => { /* to the while loop */
109 my $assign = $ASSIGN;
110 $BLOCK[1]->delete($ASSIGN);
111 $BLOCK[0]->insert_before($WHILE, $assign);
112 }
113 },
114 #outputfile => 'main.pm',
115 firstline => 104,
116);

```

Al ejecutar el programa con la entrada "a =1000; c = 1; while (a) { c = c\*a; b = 5; a = a-1 }\n" obtenemos el árbol modificado:

```

bless({ 'children' => [
 bless({ 'children' => [# a = 1000
 | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL'),
 | bless({ 'children' => [
 | bless({ 'children' => [], 'attr' => '1000', 'token' => 'NUM' }, 'TERMINAL')
 |]
 | }, 'NUM')
 |]
}, 'ASSIGN'),
bless({ 'children' => [# c = 1
 | bless({ 'children' => [], 'attr' => 'c', 'token' => 'VAR' }, 'TERMINAL'),
 | bless({ 'children' => [bless({ 'children' => [], 'attr' => '1', 'token' => 'NUM' },
 |]
 | }, 'NUM')
 |]
}, 'ASSIGN'),
bless({ 'children' => [# b = 5 moved out of loop
 | bless({ 'children' => [], 'attr' => 'b', 'token' => 'VAR' }, 'TERMINAL'),
 | bless({ 'children' => [bless({ 'children' => [], 'attr' => '5', 'token' => 'NUM' },
 |]
 | }, 'NUM')
 |]
}, 'ASSIGN'),
bless({ 'children' => [# while

```

```

| bless({ 'children' => [# (a)
| bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL')
|]
| }, 'VAR'),
| bless({ 'children' => [# BLOCK {}
| | bless({ 'children' => [# c = c * a
| | | bless({ 'children' => [], 'attr' => 'c', 'token' => 'VAR' }, 'TERMINAL'),
| | | bless({ 'children' => [
| | | | bless({ 'children' => [
| | | | | bless({ 'children' => [], 'attr' => 'c', 'token' => 'VAR' }, 'TERMINAL')
| | | |]
| | | | }, 'VAR'),
| | | | bless({ 'children' => [
| | | | | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL')
| | | |]
| | | | }, 'VAR')
| | |]
| | | }, 'TIMES')
| |]
| }, 'ASSIGN'),
| | bless({ 'children' => [# a = a - 1
| | | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL'),
| | | bless({ 'children' => [
| | | | bless({ 'children' => [
| | | | | bless({ 'children' => [], 'attr' => 'a', 'token' => 'VAR' }, 'TERMINAL')
| | | |]
| | | | }, 'VAR'),
| | | | bless({ 'children' => [
| | | | | bless({ 'children' => [], 'attr' => '1', 'token' => 'NUM' }, 'TERMINAL')
| | | |]
| | | | }, 'NUM')
| | |]
| | | }, 'MINUS')
| |]
| }, 'ASSIGN')
|]
| }, 'BLOCK')
|]
}, 'WHILE')
]
}, 'BLOCK');

```

### Expresión regular árbol estrella

Una *expresión regular árbol estrella* casa con la secuencia mas corta de hijos del nodos tal que la siguiente expresión regular árbol casa. Si no existe expresión regular árbol siguiente - esto es, la expresión regular estrella es la última de la lista como en  $A(B(C, .), *)$ - la expresión regular estrella casará con todos los nodos hijo a partir del último casamiento. Una expresión regular árbol array se escribe insertando el símbolo  $*$  en la expresión regular árbol. Las listas de nodos que han casado con la expresiones regulares árbol estrella quedaran en las variables léxicas @W\_0, @W\_1, @W\_2, etc. En este sentido una expresión regular árbol estrella no es mas que una abreviación para la expresión regular árbol @W\_# siendo # el número de orden de aparición.

### Parámetros Pasados a una Subrutina de Transformación Árbol

Como se ha mencionado anteriormente el compilador de expresiones regulares árbol traduce cada transformación árbol en una subrutina Perl. Con mayor precisión: se crea un objeto `Parse::Eyapp::YATW` que es el encargado de gestionar la transformación. Para que una subrutina pueda ser convertida en un objeto YATW deben ajustarse al *Protocolo YATW de LLamada*. Actualmente (2006) la subrutina asociada con un objeto YATW es llamada como sigue:

```
pattern_sub(
 $_[0], # Node being visited
 $_[1], # Father of this node
 $index, # Index of this node in @Father->children
 $self, # The YATW pattern object
);
```

Los cuatro argumentos tienen el siguiente significado:

1. El nódo del árbol que esta siendo visitado
2. El padre de dicho nodo
3. El índice del nodo (`$_[0]`) en la lista de nodos del padre
4. Una referencia al objeto YATW

La subrutina debe devolver cierto (TRUE) si se produce matching y falso en otro caso. Recuerde que el método `s` de los nodos (no así el de los objetos YATW) permanecerá aplicando las transformaciones hasta que no se produzcan emparejamientos. Por tanto es importante asegurarse cuando se usa la forma `$node->s(@transformations)` que la aplicación reiterada de las transformaciones conduce a situaciones en las que eventualmente las subrutinas retornan el valor falso.

Es posible que el protocolo de llamada YATW cambie en un futuro próximo.

## 10.7. La opción SEVERITY

La opción SEVERITY del constructor `Parse::Eyapp::Treeregexp::new` controla la forma en la que se interpreta el éxito de un casamiento en lo que se refiere al número de hijos del nodo. Para ilustrar su uso consideremos el siguiente ejemplo que hace uso de `Rule6` la gramática que fue introducida en la sección ?? (página ??).

```
nereida:~/src/perl/YappWithDefaultAction/examples> cat -n numchildren.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Rule6;
4 use Data::Dumper;
5 use Parse::Eyapp::Treeregexp;
6 use Parse::Eyapp::Node;
7
8 our @all;
9 my $severity = shift || 0;
10
11 my $transform = Parse::Eyapp::Treeregexp->new(STRING => q{
12 zero_times_whatever: TIMES(NUM($x)) and { $x->{attr} == 0 } => { $_[0] = $NUM }
13 },
14 SEVERITY => $severity,
15);
16
17 $transform->generate;
```



```

18 $Data::Dumper::Indent = 1;
19 my $parser = new Rule6();
20 my $t = $parser->Run;
21 $t->s(@all);
22 print Dumper($t);

```

Este programa obtiene el nivel de severidad a usar desde la línea de comandos (línea 9). Nótese que la especificación de TIMES en la transformación `zero_times_whatever` este aparece con un único hijo. Existen varias interpretaciones de la expresión que se corresponden con los niveles de SEVERITY :

- 0: Quiero decir que tiene al menos un hijo. No me importa si tiene mas
- 1: Para que case tiene que tener exactamente un hijo
- 2: Para que case tiene que tener exactamente un hijo. Si aparece un nodo TIMES con un número de hijos diferente quiero ser avisado
- 3: Para que case tiene que tener exactamente un hijo. Si aparece un nodo TIMES con un número de hijos diferente quiero que se considere un error (mis nodos tiene aridad fija)

En la siguiente ejecución el nivel especificado es cero. La expresión `0*2` casa y es modificada.

```

nereida:~/src/perl/YappWithDefaultAction/examples> numchildren.pl 0
0*2
$VAR1 = bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => 0,
 'token' => 'NUM'
 }, 'TERMINAL')
]
}, 'NUM');

```

En la siguiente ejecución el nivel especificado es uno. La expresión `0*2` no casa pero no se avisa ni se considera un error la presencia de un nodo TIMES con una aridad distinta.

```

nereida:~/src/perl/YappWithDefaultAction/examples> numchildren.pl 1
0*2
$VAR1 = bless({
 'children' => [
 bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => '0',
 'token' => 'NUM'
 }, 'TERMINAL')
]
 }, 'NUM'),
 bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => '2',
 'token' => 'NUM'
 }, 'TERMINAL')
]
 }, 'NUM')
]
}, 'NUM');

```

```

]
 }, 'NUM')
]
}, 'TIMES');

```

En la siguiente ejecución el nivel especificado es dos. La expresión `0*2` no casa y se avisa de la presencia de un nodo `TIMES` con una aridad distinta:

```

nereida:~/src/perl/YappWithDefaultAction/examples> numchildren.pl 2
0*2
Warning! found node TIMES with 2 children.
Expected 1 children (see line 12 of numchildren.pl)"
$VAR1 = bless({
 'children' => [
 bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => '0',
 'token' => 'NUM'
 }, 'TERMINAL')
]
 }, 'NUM'),
 bless({
 'children' => [
 bless({
 'children' => [],
 'attr' => '2',
 'token' => 'NUM'
 }, 'TERMINAL')
]
 }, 'NUM')
]
}, 'TIMES');

```

En la última ejecución el nivel especificado es tres. El programa se detiene ante la presencia de un nodo `TIMES` con una aridad distinta:

```

nereida:~/src/perl/YappWithDefaultAction/examples> numchildren.pl 3
0*2
Error! found node TIMES with 2 children.
Expected 1 children (see line 12 of numchildren.pl)"
at (eval 3) line 28

```

# Capítulo 11

## Análisis Sintáctico con yacc

### 11.1. Introducción a yacc

Veamos un ejemplo sencillo de analizador sintáctico escrito en yacc. La gramática se especifica entre las dos líneas de `%%`. Por defecto, el símbolo de arranque es el primero que aparece, en este caso `list`. En bison es posible hacer que otro variable lo sea utilizando la declaración `%start`:

#### Ejemplo: La Calculadora en yacc

**Programa 11.1.1.** *Calculadora elemental. Analizador sintáctico.*

```
nereida:~/src/precedencia/hoc1> cat -n hoc1.y
 1 %{
 2 /* File: /home/pl/src/precedencia/hoc1/hoc1.y */
 3 #define YYSTYPE double
 4 #include <stdio.h>
 5 %}
 6 %token NUMBER
 7 %left '+' '-'
 8 %left '*' '/'
 9 %%
10 list
11 :
12 | list '\n'
13 | list expr { printf("%.8g\n", $2); }
14 ;
15
16 expr
17 : NUMBER { $$ = $1; }
18 | expr '+' expr { $$ = $1 + $3; }
19 | expr '-' expr { $$ = $1 - $3; }
20 | expr '*' expr { $$ = $1 * $3; }
21 | expr '/' expr { $$ = $1 / $3; }
22 ;
23
24 %%
25
26 extern FILE * yyin;
27
28 main(int argc, char **argv) {
29 if (argc > 1) yyin = fopen(argv[1], "r");
30 yydebug = 1;
```

```

31 yyparse();
32 }
33
34 yyerror(char *s) {
35 printf("%s\n",s);
36 }

```

La macro `YYSTYPE` (línea 3) contiene el tipo del valor semántico. Si no se declara se asume `int`.

El fichero `yyin` (líneas 26 y 29) es definido en el fichero conteniendo el analizador léxico `lex.yy.c`. Refiere al fichero de entrada conteniendo el texto a analizar.

Al poner la variable `yydebug` a 1 activamos el modo depuración. Para que la depuración se haga efectiva es necesario definir además la macro `YYDEBUG`.

El analizador sintáctico proveído por `yacc` se llama `yyparse` (línea 31). Por defecto su declaración es `int yyparse ()`

## El Analizador Léxico

### Programa 11.1.2. Calculadora elemental. Analizador léxico:

```

nereida:~/src/precedencia/hoc1> cat -n hoc1.1
 1 %{
 2 #include "y.tab.h"
 3 extern YYSTYPE yylval;
 4 %}
 5 number [0-9]+(\\. [0-9]+)?([eE] [+-]?[0-9]+)?
 6 %%
 7 {number} { yylval = atof(yytext); return NUMBER; }
 8 .|\n { return yytext[0];}
 9 %%
10 int yywrap() { return 1; }

```

## Compilación

Al compilar el program `yacc` con la opción `-d` se produce además del fichero `y.tab.c` conteniendo el analizador sintáctico un fichero adicional de cabecera `y.tab.h` conteniendo las definiciones de los terminales:

```

nereida:~/src/precedencia/hoc1> yacc -d -v hoc1.y
nereida:~/src/precedencia/hoc1> ls -lt | head -4
total 200
-rw-rw---- 1 pl users 2857 2007-01-18 10:26 y.output
-rw-rw---- 1 pl users 35936 2007-01-18 10:26 y.tab.c
-rw-rw---- 1 pl users 1638 2007-01-18 10:26 y.tab.h
nereida:~/src/precedencia/hoc1> sed -ne '27,48p' y.tab.h | cat -n
 1 #ifndef YYTOKENTYPE
 2 # define YYTOKENTYPE
 3 /* Put the tokens into the symbol table, so that GDB and other debuggers
 4 know about them. */
 5 enum yytokentype {
 6 NUMBER = 258
 7 };
 8 #endif
 9 /* Tokens. */
10 #define NUMBER 258
..

```

```

15 #if ! defined (YYSTYPE) && ! defined (YYSTYPE_IS_DECLARED)
16 typedef int YYSTYPE;
17 # define yystype YYSTYPE /* obsolescent; will be withdrawn */
18 # define YYSTYPE_IS_DECLARED 1
19 # define YYSTYPE_IS_TRIVIAL 1
20 #endif
21
22 extern YYSTYPE yylval;

```

La variable `yylval` (líneas 3 y 7 del listado 11.1.2) es declarada por el analizador sintáctico y usada por el analizador léxico. El analizador léxico deja en la misma el valor semántico asociado con el token actual.

## Makefile

Para compilar todo el proyecto usaremos el siguiente fichero Makefile:

### Programa 11.1.3. Calculadora elemental. Makefile:

```

> cat Makefile
hoc1: y.tab.c lex.yy.c
 gcc -DYYDEBUG=1 -g -o hoc1 y.tab.c lex.yy.c
y.tab.c y.tab.h: hoc1.y
 yacc -d -v hoc1.y
lex.yy.c: hoc1.l y.tab.h
 flex -l hoc1.l
clean:
 - rm -f y.tab.c lex.yy.c *.o core hoc1

```

## Ejecución

**Ejecución 11.1.1.** *Para saber que esta haciendo el analizador, insertamos una asignación: `yydebug = 1`; justo antes de la llamada a `yyparse()` y ejecutamos el programa resultante:*

```

$ hoc1
yydebug: state 0, reducing by rule 1 (list :)
yydebug: after reduction, shifting from state 0 to state 1
2.5+3.5+1

```

*Introducimos la expresión 2.5+3.5+1. Antes que incluso ocurra la entrada, el algoritmo LR reduce por la regla  $List \rightarrow \epsilon$ .*

```

yydebug: state 1, reading 257 (NUMBER)
yydebug: state 1, shifting to state 2
yydebug: state 2, reducing by rule 4 (expr : NUMBER)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 43 ('+')
yydebug: state 4, shifting to state 5
yydebug: state 5, reading 257 (NUMBER)
yydebug: state 5, shifting to state 2
yydebug: state 2, reducing by rule 4 (expr : NUMBER)
yydebug: after reduction, shifting from state 5 to state 6
yydebug: state 6, reducing by rule 5 (expr : expr '+' expr)

```

*Observe como la declaración de la asociatividad a izquierdas `%left '+'` se traduce en la reducción por la regla 5.*

```

yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 43 ('+')
yydebug: state 4, shifting to state 5
yydebug: state 5, reading 257 (NUMBER)
yydebug: state 5, shifting to state 2
yydebug: state 2, reducing by rule 4 (expr : NUMBER)
yydebug: after reduction, shifting from state 5 to state 6
yydebug: state 6, reducing by rule 5 (expr : expr '+' expr)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 10 ('\n')
yydebug: state 4, reducing by rule 3 (list : list expr)
7

```

La reducción por la regla  $list \rightarrow list\ expr$  produce la ejecución del `printf("%.8g\n", $2);` asociado con la regla y la salida del valor 7 que constituye el atributo de *expr*.

```

yydebug: after reduction, shifting from state 0 to state 1
yydebug: state 1, shifting to state 3
yydebug: state 3, reducing by rule 2 (list : list '\n')
yydebug: after reduction, shifting from state 0 to state 1
yydebug: state 1, reading 0 (end-of-file)
$

```

En Unix la combinación de teclas `CTRL-D` nos permite generar el final de fichero.

## 11.2. Precedencia y Asociatividad

En caso de que no existan indicaciones específicas *yacc* resuelve los conflictos que aparecen en la construcción de la tabla utilizando las siguientes reglas:

1. Un conflicto *reduce-reduce* se resuelve eligiendo la producción que se listó primero en la especificación de la gramática.
2. Un conflicto *shift-reduce* se resuelve siempre en favor del *shift*

La precedencia se utiliza para modificar estos criterios. Para ello se define:

1. La precedencia de los *tokens* es determinada según el orden de declaración. La declaración de *tokens* mediante la palabra reservada `token` no modifica la precedencia. Si lo hacen las declaraciones realizadas usando las palabras `left`, `right` y `nonassoc`. Los *tokens* declarados en la misma línea tienen igual precedencia. La precedencia es mayor cuanto mas abajo en el texto. Así, en el ejemplo que sigue, el *token* `*` tiene mayor precedencia que `+` pero la misma que `/`.
2. La precedencia de una regla  $A \rightarrow \alpha$  se define como la del terminal mas a la derecha que aparece en  $\alpha$ . En el ejemplo, la producción

$$\text{expr} : \text{expr} \text{ '+' expr}$$

tiene la precedencia del *token* `+`.

3. Para decidir en un conflicto *shift-reduce* se comparan la precedencia de la regla con la del terminal que va a ser desplazado. Si la de la regla es mayor se reduce si la del *token* es mayor, se desplaza.
4. Si en un conflicto *shift-reduce* ambos la regla y el terminal que va a ser desplazado tiene la misma precedencia *yacc* considera la asociatividad, si es asociativa a izquierdas, reduce y si es asociativa a derechas desplaza. Si no es asociativa, genera un mensaje de error.

Obsérvese que, en esta situación, la asociatividad de la regla y la del *token* han de ser por fuerza, las mismas. Ello es así, porque en *yacc* los *tokens* con la misma precedencia se declaran en la misma línea.

*Por tanto es imposible declarar dos tokens con diferente asociatividad y la misma precedencia.*

5. Es posible modificar la precedencia “natural” de una regla, calificándola con un *token* específico. para ello se escribe a la derecha de la regla **prec token**, donde **token** es un *token* con la precedencia que deseamos. Vea el uso del *token dummy* en el siguiente ejercicio.

**Programa 11.2.1.** *Este programa muestra el manejo de las reglas de precedencia.*

```
%{
#define YYSTYPE double
#include <stdio.h>
%}
%token NUMBER
%left '@'
%right '&' dummy
%%
list
 :
 | list '\n'
 | list e
 ;

e : NUMBER
 | e '&' e
 | e '@' e %prec dummy
 ;

%%
extern FILE * yyin;

main(int argc, char **argv) {
 if (argc > 1) yyin = fopen(argv[1], "r");
 yydebug = 1;
 yyparse();
}

yyerror(char *s) {
 printf("%s\n", s);
}
```

**Ejercicio 11.2.1.** *Dado el programa yacc 11.2.1 Responda a las siguientes cuestiones:*

1. *Construya las tablas SLR de acciones y gotos.*
2. *Determine el árbol construido para las frases: 4@3@5, 4&3&5, 4@3&5, 4&3@5.*
3. *¿Cuál es la asociatividad final de la regla e : e '@' e?*

**Listado 11.2.1.** *Fichero y.output:*

```
0 $accept : list $end
```

```

1 list :
2 | list '\n'
3 | list e

4 e : NUMBER
5 | e '&' e
6 | e '@' e

```

|                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |                                                                                                                                                                                                                                                                                                                                                                                                                                           |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> state 0 \$accept : . list \$end (0) list : . (1)  . reduce 1  list goto 1  state 1 \$accept : list . \$end (0) list : list . '\n' (2) list : list . e (3)  \$end accept NUMBER shift 2 '\n' shift 3 . error  e goto 4  state 2 e : NUMBER . (4)  . reduce 4  state 3 list : list '\n' . (2)  . reduce 2  state 4 list : list e . (3) e : e . '&amp;' e (5) e : e . '@' e (6)  '@' shift 5 '&amp;' shift 6 \$end reduce 3 NUMBER reduce 3 '\n' reduce 3 </pre> | <pre> state 5 e : e '@' . e (6)  NUMBER shift 2 . error  e goto 7  state 6 e : e '&amp;' . e (5)  NUMBER shift 2 . error  e goto 8  state 7 e : e . '&amp;' e (5) e : e . '@' e (6) e : e '@' e . (6)  '&amp;' shift 6 \$end reduce 6 NUMBER reduce 6 '@' reduce 6 '\n' reduce 6  state 8 e : e . '&amp;' e (5) e : e '&amp;' e . (5) e : e . '@' e (6)  '&amp;' shift 6 \$end reduce 5 NUMBER reduce 5 '@' reduce 5 '\n' reduce 5 </pre> |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|



7 terminals, 3 nonterminals

7 grammar rules, 9 states

**Ejemplo 11.2.1.** *Contrasta tu respuesta con la traza seguida por el programa anterior ante la entrada 1@2&3, al establecer la variable yydebug = 1 y definir la macro YYDEBUG:*

**Ejecución 11.2.1.** \$ hocprec

```
yydebug: state 0, reducing by rule 1 (list :)
yydebug: after reduction, shifting from state 0 to state 1
1@2&3
yydebug: state 1, reading 257 (NUMBER)
yydebug: state 1, shifting to state 2
yydebug: state 2, reducing by rule 4 (e : NUMBER)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 64 ('@')
yydebug: state 4, shifting to state 5
yydebug: state 5, reading 257 (NUMBER)
yydebug: state 5, shifting to state 2
yydebug: state 2, reducing by rule 4 (e : NUMBER)
yydebug: after reduction, shifting from state 5 to state 7
yydebug: state 7, reading 38 ('&')
yydebug: state 7, shifting to state 6
```

*¿Por que se desplaza? ¿No va eso en contra de la declaración %left '@'?. ¿O quizá es porque la precedencia de @ es menor que la de &? La respuesta es que la precedencia asignada por la declaración*

$$e : e '@' e \%prec dummy$$

*cambio la asociatividad de la regla. Ahora la regla se “enfrenta” a un token, & con su misma precedencia. Al pasar a ser asociativa a derechas (debido a que dummy lo es), se debe desplazar y no reducir.*

**Ejemplo 11.2.2.** *Otra ejecución, esta vez con entrada 1&2@3. Compara tus predicciones con los resultados.*

**Ejecución 11.2.2.** \$ hocprec

```
yydebug: state 0, reducing by rule 1 (list :)
yydebug: after reduction, shifting from state 0 to state 1
1&2@3
yydebug: state 1, reading 257 (NUMBER)
yydebug: state 1, shifting to state 2
yydebug: state 2, reducing by rule 4 (e : NUMBER)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 38 ('&')
yydebug: state 4, shifting to state 6
yydebug: state 6, reading 257 (NUMBER)
yydebug: state 6, shifting to state 2
yydebug: state 2, reducing by rule 4 (e : NUMBER)
yydebug: after reduction, shifting from state 6 to state 8
yydebug: state 8, reading 64 ('@')
yydebug: state 8, reducing by rule 5 (e : e '&' e)
```

*En este caso se comparan la producción :*

$$e \rightarrow e \& e$$

*con el token @. La regla tiene mayor precedencia que el token, dado que la precedencia de la regla es la de &.*

### 11.3. Uso de union y type

En general, los atributos asociados con las diferentes variables sintácticas y terminales tendrán tipos de datos distintos. Para ello, `yacc` provee la declaración `%union`.

La declaración `%union` especifica la colección de posibles tipos de datos de `yyval` y de los atributos `$1`, `$2`, ...

He aquí un ejemplo:

```
%union {
 double val;
 int index;
}
```

Esto dice que los dos tipos de alternativas son `double` y `int`. Se les han dado los nombres `val` y `index`.

```
%token <val> NUMBER
%token <index> VAR
%type <val> expr
%right '='
%left '+' '-'
%left '*' '/'
```

Estos nombres `<val>` e `<index>` se utilizan en las declaraciones de `%token` y `%type` para definir el tipo del correspondiente atributo asociado.

Dentro de las acciones, se puede especificar el tipo de un símbolo insertando `<tipo>` después del `$` que referencia al atributo. En el ejemplo anterior podríamos escribir `$<val>1` para indicar que manipulamos el atributo del primer símbolo de la parte derecha de la regla como si fuera un `double`.

La información que provee la declaración `%union` es utilizada por `yacc` para realizar la sustitución de las referencias textuales/formales (`$$`, `$1`, ... `$| $\alpha$ |`) a los atributos de los símbolos que conforman la regla (que pueden verse como parámetros formales de las acciones) por las referencias a las zonas de memoria en las que se guardan (que están asociadas con los correspondientes estados de la pila) cuando tiene lugar la reducción en el algoritmo de análisis LR:

```
case "reduce A \rightarrow α " :
 execute("reduce A \rightarrow α ", top(| α |-1), ... , top(0));
 pop(| α |);
 push(goto[top(0)][A]);
 break;
```

Así, `yacc` es capaz de insertar el código de ahorrado de tipos correcto. Ello puede hacerse porque se conocen los tipos asociados con los símbolos en la parte derecha de una regla, ya que han sido proveídos en las declaraciones `%union`, `%token` y `%type`.

### 11.4. Acciones en medio de una regla

A veces necesitamos insertar una acción en medio de una regla. Una acción en medio de una regla puede hacer referencia a los atributos de los símbolos que la preceden (usando `$n`), pero no a los que le siguen.

Cuando se inserta una acción `{action1}` para su ejecución en medio de una regla  $A \rightarrow \alpha\beta$ :

$$A \rightarrow \alpha \{action_1\} \beta \{action_2\}$$

`yacc` crea una variable sintáctica temporal  $T$  e introduce una nueva regla:

1.  $A \rightarrow \alpha T \beta \{action_2\}$

## 2. $T \rightarrow \epsilon \{action_1\}$

Las acciones en mitad de una regla cuentan como un símbolo mas en la parte derecha de la regla. Asi pues, en una acción posterior en la regla, se deberán referenciar los atributos de los símbolos, teniendo en cuenta este hecho.

Las acciones en mitad de la regla pueden tener un atributo. La acción en cuestión puede hacer referencia a ese atributo mediante \$\$, y las acciones posteriores en la regla se referirán a él como \$n, siendo n su número de orden en la parte derecha. Dado que no existe un símbolo explícito que identifique a la acción, no hay manera de que el programador declare su tipo. Sin embargo, es posible utilizar la construcción \$<valtipo># para especificar la forma en la que queremos manipular su atributo.

Na hay forma de darle, en una acción a media regla, un valor al atributo asociado con la variable en la izquierda de la regla de producción (ya que \$\$ se refiere al atributo de la variable temporal utilizada para introducir la acción a media regla).

**Programa 11.4.1.** *El siguiente programa ilustra el uso de %union y de las acciones en medio de una regla.*

```
%{
#include <string.h>
char buffer[256];
#define YYDEBUG 1
%}
%union {
 char tA;
 char *tx;
}
%token <tA> A
%type <tx> x
%%
s : x { *$1 = '\0'; printf("%s\n",buffer); } '\n' s
 |
 ;

x : A { $$ = buffer + sprintf(buffer,"%c",$1); }
 | A { $<tx>$ = strdup("**"); } x
 { $$ = $3 + sprintf($3,"%s%c",$<tx>2,$1); free($2); }
 ;

%%

main() {
 yydebug=1;
 yyparse();
}

yyerror(char *s) {
 printf("%s\n",s);
}
```

**Programa 11.4.2.** *El analizador léxico utilizado es el siguiente:*

```
%{
#include "y.tab.h"
%}
%%
```

```
[\t]+
[a-zA-Z0-9] { yylval.tA = yytext[0]; return A; }
(.|\n) { return yytext[0]; }
%%
yywrap() { return 1; }
```

**Ejemplo 11.4.1.** Considere el programa yacc 11.4.1. ¿Cuál es la salida para la entrada ABC?

La gramática inicial se ve aumentada con dos nuevas variables sintácticas temporales y dos reglas  $t_1 \rightarrow \epsilon$  y  $t_2 \rightarrow \epsilon$ . Además las reglas correspondientes pasan a ser:  $s \rightarrow xt_1s$  y  $x \rightarrow At_2x$ . El análisis de la entrada ABC nos produce el siguiente árbol anotado:

**Ejecución 11.4.1.** Observe la salida de la ejecución del programa 11.4.1. La variable “temporal” creada por yacc para la acción en medio de la regla

$$s \rightarrow x \{ *\$1 = '\0'; printf("%s\n",buffer); \} '\n' s$$

se denota por  $\$ \$1$ . La asociada con la acción en medio de la regla

$$x \rightarrow A \{ \$<tx>\$ = strdup("**"); \} x$$

se denota  $\$ \$2$ .

```
$ yacc -d -v media4.y ; flex -l medial.l ; gcc -g y.tab.c lex.yy.c
$ a.out
ABC
yydebug: state 0, reading 257 (A)
yydebug: state 0, shifting to state 1
yydebug: state 1, reading 257 (A)
yydebug: state 1, reducing by rule 5 ($$2 :)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, shifting to state 1
yydebug: state 1, reading 257 (A)
yydebug: state 1, reducing by rule 5 ($$2 :)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, shifting to state 1
yydebug: state 1, reading 10 ('\n')
yydebug: state 1, reducing by rule 4 (x : A)
yydebug: after reduction, shifting from state 4 to state 6
yydebug: state 6, reducing by rule 6 (x : A $$2 x)
yydebug: after reduction, shifting from state 4 to state 6
yydebug: state 6, reducing by rule 6 (x : A $$2 x)
yydebug: after reduction, shifting from state 0 to state 3
yydebug: state 3, reducing by rule 1 ($$1 :)
C**B**A
yydebug: after reduction, shifting from state 3 to state 5
yydebug: state 5, shifting to state 7
yydebug: state 7, reading 0 (end-of-file)
yydebug: state 7, reducing by rule 3 (s :)
yydebug: after reduction, shifting from state 7 to state 8
yydebug: state 8, reducing by rule 2 (s : x $$1 '\n' s)
yydebug: after reduction, shifting from state 0 to state 2
```

**Ejemplo 11.4.2.** ¿Que ocurre si en el programa 11.4.1 adelantamos la acción intermedia en la regla

$$x \rightarrow A \{ \$<tx>\$ = strdup("**"); \} x$$

y la reescribimos

$$x \rightarrow \{ \$\langle tx \rangle \$ = \text{strdup}("***"); \} \text{ A } x?$$

**Ejecución 11.4.2.** *En tal caso obtendremos:*

```
$ yacc -d -v media3.y
yacc: 1 rule never reduced
yacc: 3 shift/reduce conflicts.
```

*¿Cuáles son esos 3 conflictos?*

**Listado 11.4.1.** *El fichero y.output comienza enumerando las reglas de la gramática extendida:*

```
1 0 $accept : s $end
2
3 1 $$1 :
4
5 2 s : x $$1 '\n' s
6 3 |
7
8 4 x : A
9
10 5 $$2 :
11
12 6 x : $$2 A x
13 ^L
```

*A continuación nos informa de un conflicto en el estado 0. Ante el token A no se sabe si se debe desplazar al estado 1 o reducir por la regla 5: `$$2 : .`*

```
14 0: shift/reduce conflict (shift 1, reduce 5) on A
15 state 0
16 $accept : . s $end (0)
17 s : . (3)
18 $$2 : . (5)
19
20 A shift 1
21 $end reduce 3
22
23 s goto 2
24 x goto 3
25 $$2 goto 4
```

*Observe que, efectivamente, `$$2 : .` esta en la clausura del estado de arranque del NFA (`$accept : . s $end`). Esto es así, ya que al estar el marcador junto a `x`, estará el ítem `s : . x $$1 '\n' s` y de aquí que también este `x : . $$2 A x`.*

*Además el token A está en el conjunto FOLLOW(`$$2`) (Basta con mirar la regla 6 para confirmarlo). Por razones análogas también está en la clausura del estado de arranque el ítem `x : . A` que es el que motiva el desplazamiento al estado 1.*

*La dificultad para yacc se resolvería si dispusiera de información acerca de cual es el token que viene después de la A que causa el conflicto.*

**Ejercicio 11.4.1.** *¿Que acción debe tomarse en el conflicto del ejemplo 11.4.2 si el token que viene después de A es `\n`? ¿Y si el token es A? ¿Se debe reducir o desplazar?*

## 11.5. Recuperación de Errores

Las entradas de un traductor pueden contener errores. El lenguaje `yacc` proporciona un *token* especial, `error`, que puede ser utilizado en el programa fuente para extender el traductor con producciones de error que lo doten de cierta capacidad para recuperarse de una entrada errónea y poder continuar analizando el resto de la entrada.

**Ejecución 11.5.1.** *Consideremos lo que ocurre al ejecutar el programa `yacc` 11.1.1 con una entrada errónea:*

```
$ hoc1
yydebug: state 0, reducing by rule 1 (list :)
yydebug: after reduction, shifting from state 0 to state 1
3--2
yydebug: state 1, reading 257 (NUMBER)
yydebug: state 1, shifting to state 2
yydebug: state 2, reducing by rule 4 (expr : NUMBER)
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 45 (illegal-symbol)
syntax error
yydebug: error recovery discarding state 4
yydebug: error recovery discarding state 1
yydebug: error recovery discarding state 0
```

Después de detectar el mensaje `yacc` emite el mensaje `syntax error` y comienza a sacar estados de la pila hasta que esta se vacía.

**Programa 11.5.1.** *La conducta anterior puede modificarse si se introducen “reglas de recuperación de errores” como en la siguiente modificación del programa 11.1.1:*

```
%{
#define YYSTYPE double
#define YYDEBUG 1
#include <stdio.h>
%}
%token NUMBER
%left '-'
%%
list
:
| list '\n'
| list error '\n' { yyerrok; }
| list expr { printf("%.8g\n", $2); }
;

expr
: NUMBER { $$ = $1; }
| expr '-' expr { $$ = $1 - $3; }
;

%%
```

La regla `list → list error '\n'` es una producción de error. La idea general de uso es que, a través de la misma, el programador le indica a `yacc` que, cuando se produce un error dentro de una expresión, descarte todos los *tokens* hasta llegar al retorno del carro y prosiga con el análisis. Además, mediante la llamada a la macro `yyerrok` el programador anuncia que, si se alcanza este punto, la

recuperación puede considerarse “completa” y que yacc puede emitir a partir de ese momento mensajes de error con la seguridad de que no son consecuencia de un comportamiento inestable provocado por el primer error.

**Algoritmo 11.5.1.** *El esquema general del algoritmo de recuperación de errores usado por la versión actual de yacc es el siguiente:*

1. Cuando se encuentra ante una acción de error, el analizador genera un token **error**.
2. A continuación pasa a retirar estados de la pila hasta que descubre un estado capaz de transitar ante el token **error**.
3. En este punto transita al estado correspondiente a desplazar el token **error**.
4. Entonces lee tokens y los descarta hasta encontrar uno que sea aceptable.
5. Sólo se envían nuevos mensajes de error una vez asimilados (desplazados) tres símbolos terminales. De este modo se intenta evitar la aparición masiva de mensajes de error.

**Algoritmo 11.5.2.** *El cuerpo principal del analizador LR permanece sin demasiados cambios:*

```
goodtoken = 3; b = yylex();
for(; ;) {
 s = top(); a = b;
 switch (action[s][a]) {
 case "shift t" : push(t); b = yylex(); goodtoken++; break;
 case "reduce A -> alpha" :
 pop(strlen(alpha));
 push(goto[top()][A]);
 break;
 case "accept" : return (1);
 default : if (errorrecovery("syntax error")) return (ERROR);
 }
}
```

**Algoritmo 11.5.3.** *El siguiente pseudocódigo es una reescritura mas detallada del algoritmo 11.5.1. Asumimos que las funciones pop() y popstate() comprueban que hay suficientes elementos en la pila para retirar. En caso contrario se emitirá un mensaje de error y se terminará el análisis.*

```
errorrecovery(char * s) {
 if (goodtoken > 2) {
 yyerror(s); goodtoken = 0;
 }
 while (action[s][error] != shift)
 popstate(s);
 push(goto[s][error]);
 s = top();
 while (action[s][a] == reduce A -> alpha) {
 pop(strlen(alpha));
 push(goto[top()][A]);
 s = top();
 }
 switch (action[s][a]) {
 case "shift t" :
 push(t);
 b = yylex();
```

```

 goodtoken++;
 RETURN RECOVERING;
case "accept" : return (ERROR);
default :
 do b = yylex();
 while ((b != EOF)&&(action[s][b] == error));
 if (b == EOF)
 return (ERROR);
 else
 RETURN RECOVERING
}

```

Parecen existir diferencias en la forma en la que *bison* y *yacc* se recuperan de los errores.

#### Ejecución 11.5.2. Ejecutemos el programa 11.5.1 con yacc.

```

$ yacc -d hoc1.y; flex -l hoc1.l; gcc y.tab.c lex.yy.c; a.out
yydebug: state 0, reducing by rule 1 (list :)
yydebug: after reduction, shifting from state 0 to state 1
2--3-1
yydebug: state 1, reading 257 (NUMBER)
yydebug: state 1, shifting to state 3
yydebug: state 3, reducing by rule 5 (expr : NUMBER)
yydebug: after reduction, shifting from state 1 to state 5
yydebug: state 5, reading 45 ('-')
yydebug: state 5, shifting to state 7
yydebug: state 7, reading 45 ('-')
syntax error

```

Puesto que es el primer error, se cumple que (`goodtoken > 2`), emitiéndose el mensaje de error. Ahora comienzan a ejecutarse las líneas:

```

 while (!(action[s][error] != shift)) popstate(s);

```

que descartan los estados, hasta encontrar el estado que contiene el ítem

```

list → list ↑ error '\n'

```

```

yydebug: error recovery discarding state 7
yydebug: error recovery discarding state 5
yydebug: state 1, error recovery shifting to state 2

```

Una vez en ese estado, transitamos con el token `error`,

```

yydebug: state 2, error recovery discards token 45 ('-')
yydebug: state 2, reading 257 (NUMBER)
yydebug: state 2, error recovery discards token 257 (NUMBER)
yydebug: state 2, reading 45 ('-')
yydebug: state 2, error recovery discards token 45 ('-')
yydebug: state 2, reading 257 (NUMBER)
yydebug: state 2, error recovery discards token 257 (NUMBER)

```

Se ha procedido a descartar tokens hasta encontrar el retorno de carro, ejecutando las líneas:

```

 b = yylex(); while ((b != EOF)&&(action[s][b] == error));

```

```

yydebug: state 2, reading 10 ('\n')
yydebug: state 2, shifting to state 6
yydebug: state 6, reducing by rule 3 (list : list error '\n')
yydebug: after reduction, shifting from state 0 to state 1

```

Al reducir por la regla de error, se ejecuta `yyerrok` y *yacc* reestablece el valor de `goodtoken`. Si se producen nuevos errores serán señalados.



## 11.6. Recuperación de Errores en Listas

Aunque no existe un método exacto para decidir como ubicar las reglas de recuperación de errores, en general, los símbolos de error deben ubicarse intentado satisfacer las siguientes reglas:

- Tan cerca como sea posible del símbolo de arranque.
- Tan cerca como sea posible de los símbolos terminales.
- Sin introducir nuevos conflictos.

**Esquema 11.6.1.** *En el caso particular de las listas, se recomienda seguir el siguiente esquema:*

| <i>Construct</i>         | <i>EBNF</i>   | <i>yacc input</i>                                                                      |
|--------------------------|---------------|----------------------------------------------------------------------------------------|
| <i>optional sequence</i> | $x : \{y\}$   | $x : /* \text{null} */$<br>$xyyyerrok; $<br>$xerror $                                  |
| <i>sequence</i>          | $x : y\{y\}$  | $x : y$<br>$xyyyerrok; $<br>$error $<br>$xerror $                                      |
| <i>list</i>              | $x : y\{Ty\}$ | $x : y$<br>$xTyyyerrok; $<br>$error $<br>$xerror $<br>$xerroryyyerrok; $<br>$xTerror $ |

**Programa 11.6.1.** *Para comprobar el funcionamiento y la validez de la metodología esquematizada en el esquema 11.6.1, consideremos los contenidos del fichero `error.y`. En el se muestra el tercer caso  $x:y\{Ty\}$  con  $x = \text{list}$ ,  $T = ,$  e  $y = \text{NUMBER}$ :*

```
%{
#include <stdio.h>
void put(double x);
void err(int code);
%}

%union {
 double val;
}
%token <val>NUMBER
%%
command
:
| command list '\n' { yyerrok; }
;

list
: NUMBER { put($1); }
| list ',' NUMBER { put($3); yyerrok; }
| error { err(1); }
| list error { err(2); }
| list error NUMBER { err(3); put($3); yyerrok; }
| list ',' error { err(4); }
;
```

```

%%
void put(double x) {
 printf("%.21lf\n",x);
}

void err(int code) {
 printf("err %d\n",code);
}

main() {
 yydebug = 1;
 yyparse();
}

yyerror(char *s) {
 printf("%s\n",s);
}

```

**Listado 11.6.1.** *La compilación con yacc da lugar a una tabla ligeramente diferente de la producida por bison. El fichero y.output contiene la tabla:*

```

0 $accept : command $end

1 command :
2 | command list '\n'

3 list : NUMBER
4 | list ',' NUMBER
5 | error
6 | list error
7 | list error NUMBER
8 | list ',' error

state 0
$accept : . command $end (0)
command : . (1)

. reduce 1

command goto 1

state 1
$accept : command . $end (0)
command : command . list '\n' (2)

$end accept
error shift 2
NUMBER shift 3
. error

list goto 4

```

```

state 2
list : error . (5)

. reduce 5

state 3
list : NUMBER . (3)

. reduce 3

state 4
command : command list . '\n' (2)
list : list . ',' NUMBER (4)
list : list . error (6)
list : list . error NUMBER (7)
list : list . ',' error (8)

error shift 5
'\n' shift 6
',' shift 7
. error

state 5
list : list error . (6)
list : list error . NUMBER (7)

NUMBER shift 8
error reduce 6
'\n' reduce 6
',' reduce 6

state 6
command : command list '\n' . (2)

. reduce 2

state 7
list : list ',' . NUMBER (4)
list : list ',' . error (8)

error shift 9
NUMBER shift 10
. error

state 8
list : list error NUMBER . (7)

```

. reduce 7

state 9

list : list ',' error . (8)

. reduce 8

state 10

list : list ',' NUMBER . (4)

. reduce 4

5 terminals, 3 nonterminals

9 grammar rules, 11 states

**Ejecución 11.6.1.** *La ejecución del programa generado por yacc es como sigue:*

> error

yydebug: state 0, reducing by rule 1 (command :)

yydebug: after reduction, shifting from state 0 to state 1

10 20

yydebug: state 1, reading 257 (NUMBER)

yydebug: state 1, shifting to state 3

yydebug: state 3, reducing by rule 3 (list : NUMBER)

10.0

yydebug: after reduction, shifting from state 1 to state 4

yydebug: state 4, reading 257 (NUMBER)

syntax error

yydebug: state 4, error recovery shifting to state 5

yydebug: state 5, shifting to state 8

yydebug: state 8, reducing by rule 7 (list : list error NUMBER)

err 3

20.0

yydebug: after reduction, shifting from state 1 to state 4

yydebug: state 4, reading 10 ('\n')

yydebug: state 4, shifting to state 6

yydebug: state 6, reducing by rule 2 (command : command list '\n')

yydebug: after reduction, shifting from state 0 to state 1

10;20 30

yydebug: state 1, reading 257 (NUMBER)

yydebug: state 1, shifting to state 3

yydebug: state 3, reducing by rule 3 (list : NUMBER)

10.0

yydebug: after reduction, shifting from state 1 to state 4

yydebug: state 4, reading 59 (illegal-symbol)

syntax error

yydebug: state 4, error recovery shifting to state 5

yydebug: state 5, error recovery discards token 59 (illegal-symbol)

yydebug: state 5, reading 257 (NUMBER)

yydebug: state 5, shifting to state 8

yydebug: state 8, reducing by rule 7 (list : list error NUMBER)

```

err 3
20.0
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 257 (NUMBER)
syntax error
yydebug: state 4, error recovery shifting to state 5
yydebug: state 5, shifting to state 8
yydebug: state 8, reducing by rule 7 (list : list error NUMBER)
err 3
30.0
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 10 ('\n')
yydebug: state 4, shifting to state 6
yydebug: state 6, reducing by rule 2 (command : command list '\n')
yydebug: after reduction, shifting from state 0 to state 1
3,
yydebug: state 1, reading 257 (NUMBER)
yydebug: state 1, shifting to state 3
yydebug: state 3, reducing by rule 3 (list : NUMBER)
3.0
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, reading 44 (',')
yydebug: state 4, shifting to state 7
yydebug: state 7, reading 10 ('\n')
syntax error
yydebug: state 7, error recovery shifting to state 9
yydebug: state 9, reducing by rule 8 (list : list ', ' error)
err 4
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, shifting to state 6
yydebug: state 6, reducing by rule 2 (command : command list '\n')
yydebug: after reduction, shifting from state 0 to state 1
#
yydebug: state 1, reading 35 (illegal-symbol)
syntax error
yydebug: state 1, error recovery shifting to state 2
yydebug: state 2, reducing by rule 5 (list : error)
err 1
yydebug: after reduction, shifting from state 1 to state 4
yydebug: state 4, error recovery discards token 35 (illegal-symbol)
yydebug: state 4, reading 10 ('\n')
yydebug: state 4, shifting to state 6
yydebug: state 6, reducing by rule 2 (command : command list '\n')
yydebug: after reduction, shifting from state 0 to state 1
yydebug: state 1, reading 0 (end-of-file)

```

## Capítulo 12

# Análisis de Ámbito

### 12.1. Análisis de Ámbito: Conceptos

#### El Problema del Análisis de Ámbito

En los lenguajes de programación *name binding* (binding = encuadernado, encarpado, ligadura, unificación) es la asociación de valores con identificadores. Decimos de un identificador ligado a un valor que es una *referencia* a dicho valor. El concepto de binding es un concepto proveído por los lenguajes de programación: a nivel de máquina no existe el concepto de binding, de relación (nombre, valor). El concepto de Binding esta intimamente relacionado con el concepto de ámbito *scope*), ya que el *análisis de ámbito* es la determinación de las relaciones de binding.

El problema del análisis de ámbito sería sencillo sino fuera porque los lenguajes de programación suelen permitir el uso del mismo *nombre*<sup>1</sup> para denotar distintos elementos de un programa. Es por ello que es necesario determinar que definición o declaración se aplica a una determinada ocurrencia de un elemento.

**Definición 12.1.1.** *En un lenguaje de programación, una declaración es un constructo sintáctico que define y provee información sobre un nombre. La declaración provee información sobre las propiedades asociadas con el uso del nombre: 'este nombre es una función que recibe enteros y retorna enteros', 'este nombre puede ser usado para referirse a listas de enteros y es visible sólo en el ámbito actual', etc.*

**Definición 12.1.2.** *Las reglas de ámbito de un lenguaje determinan que declaración del nombre es la que se aplica cuando el nombre es usado.*

#### Binding Estático y Binding Dinámico

En la definición anterior no se especifica en que momento se resuelve la correspondencia (nombre, definición).

Se habla de *static binding* cuando las reglas y la resolución de la correspondencia (nombre, definición) puede ser resuelta en tiempo de compilación, a partir del análisis del texto del programa fuente (también se habla de *early binding*).

Por el contrario cuando se habla de *dynamic binding* cuando la determinación de que definición se aplica a un nombre es establecida en tiempo de ejecución (también se denomina *late binding* o *virtual binding*).

Un ejemplo de static binding es una llamada a una función en C: la función referenciada por un identificador no puede cambiarse en tiempo de ejecución. Un ejemplo de binding dinámico puede ocurrir cuando se trabaja con métodos polimorfos en un lenguaje de programación orientada a objetos, ya que la definición completa del tipo del objeto no se conoce hasta el momento de la ejecución.

El siguiente ejemplo de Dynamic binding tomado de la wikipedia, ilustra el binding dinámico.

Supongamos que todas las formas de vida son mortales. En OOP podemos decir que la clase **Persona** y la clase **Planta** deben implementar los métodos de **Mortal**, el cual contiene el método **muere**.

---

<sup>1</sup>Usamos el término nombre y no identificador ya que este último tiene una connotación mas precisa

Las personas y las plantas mueren de forma diferente, por ejemplo las plantas no tienen un corazón que se detenga. Dynamic binding es la práctica de determinar que definición/declaración se aplica a un método en tiempo de ejecución:

```
void mata(Mortal m) {
 m.muere();
}
```

No está claro cuál es la clase actual de `m`, una persona o una planta. Ambos `Planta.muere` y `Persona.muere` pueden ser invocados. Cuando se usa dynamic binding, el objeto `m` es examinado en tiempo de ejecución para determinar qué método es invocado. Esto supone una 'renuncia' por parte del lenguaje y su compilador a obtener una definición completa del objeto.

**Definición 12.1.3.** *Cuando se usa static binding, la parte del texto del programa al cual se aplica la declaración de un nombre se denomina ámbito de la declaración*

**Ejercicio 12.1.1.** *En el siguiente código existen dos definiciones para el nombre `one`, una en la línea 11 y otra en la línea 20.*

```
pl@europa:~/src/perl/perltesting$ cat -n ds.pl
 1 package Father;
 2 use warnings;
 3 use strict;
 4
 5 sub new {
 6 my $class = shift;
 7
 8 bless { @_ }, $class;
 9 }
10
11 sub one {
12 "Executing Father::one\n";
13 }
14
15 package Child;
16 use warnings;
17 use strict;
18 our @ISA = 'Father';
19
20 sub one {
21 "Executing Child::one\n";
22 }
23
24 package main;
25
26 for (1..10) {
27 my $class = int(rand(2)) ? 'Child' : 'Father';
28 my $c = $class->new;
29 print $c->one;
30 }
```

*¿Qué definiciones se aplican a los 10 usos del nombre `one` en la línea 28? ¿Estos usos constituyen un ejemplo de binding estático o dinámico? ¿Cuál es el ámbito de las declaraciones de `one`?*

Incluso en los casos en los que la resolución del binding se deja para el momento de la ejecución el compilador debe tener información suficiente para poder generar código. En el caso anterior, el

compilador de Perl infiere de la presencia de la flecha en `$c->one` que `one` es el nombre de una subrutina.

**Ejercicio 12.1.2.** *En el siguiente ejemplo se usa una referencia simbólica para acceder a una función:*

```
pl@europa:~/src/perl/testing$ cat -n symbolic.pl
1 use warnings;
2 use strict;
3
4 sub one {
5 "1\n";
6 }
7
8 sub two {
9 "2\n";
10 }
11
12 my $x = <>;
13 chomp($x);
14
15 no strict 'refs';
16 print &$x();
```

*Al ejecutarlo con entrada `one` obtenemos:*

```
pl@europa:~/src/perl/testing$ perl symbolic.pl
one
1
```

*¿El uso de la línea 16 es un ejemplo de binding estático o dinámico? ¿Cuál es el binding de las declaraciones de `x`?*

**Ejercicio 12.1.3.** *En el siguiente ejemplo la clase `Two` hereda de `One`.*

```
pl@europa:~/src/perl/testing$ cat -n latebinding.pl
1 package One;
2 use warnings;
3 use strict;
4
5 our $x = 1;
6 sub tutu {
7 "Inside tutu: x = $x\n";
8 }
9
10 package Two;
11 use warnings;
12 use strict;
13 our @ISA = 'One';
14
15 our $x = 2;
16
17 print Two->tutu();
```

*¿Qué definición de `$x` se aplica al uso en la línea 7? ¿Cuál será la salida del programa?*



**Definición 12.1.4.** *La tarea de asignar las ocurrencias de las declaraciones de nombres a las ocurrencias de uso de los nombres de acuerdo a las reglas de ámbito del lenguaje se denomina identificación de los nombres*

**Definición 12.1.5.** *Una ocurrencia de un nombre se dice local si está en el ámbito de una declaración que no se aplica desde el comienzo de la declaración hasta el final del texto del programa. Tal declaración es una declaración local.*

**Definición 12.1.6.** *Si, por el contrario, una ocurrencia de un nombre está en el ámbito de una declaración que se aplica desde el comienzo de la declaración hasta el final del texto del programa se dice que la declaración es una declaración global.*

**Definición 12.1.7.** *Aunque la definición anterior establece el atributo ámbito como un atributo de la declaración es usual y conveniente hablar del*

”ámbito del nombre  $x$ ”

como una abreviación de

”el ámbito de la declaración del nombre  $x$  que se aplica a esta ocurrencia de  $x$ ”

### Intervención del Programador en Tiempo de Compilación

En algunos lenguajes - especialmente en los lenguajes dinámicos- la diferenciación entre tiempo de compilación y tiempo de ejecución puede ser difusa. En el siguiente fragmento de código Perl se usa el módulo `Contextual::Return` para crear una variable cuya definición cambia con la forma de uso.

```
lhp@nereida:~/Lperl/src/testing$ cat -n context1.pl
1 #!/usr/local/bin/perl -w
2 use strict;
3 use Contextual::Return;
4
5 my $x = BOOL { 0 } NUM { 3.14 } STR { "pi" };
6
7 unless ($x) { warn "¡El famoso número $x (".(0+$x).") pasa a ser falso!\n" } # executed!

lhp@nereida:~/Lperl/src/testing$ context1.pl
¡El famoso número pi (3.14) pasa a ser falso!
```

Obsérvese que el binding de  $x$  es estático y que a los tres usos de  $x$  en la línea 7 se les asigna la definición en la línea 5. La declaración de  $x$  ocurre en lo que Perl denomina ‘tiempo de compilación’, sin embargo, el hecho de que un módulo cargado en tiempo de compilación puede ejecutar sentencias permite a `Contextual::Return` expandir el lenguaje de las declaraciones Perl.

**Ejercicio 12.1.4.** *Considere el siguiente código Perl:*

```
pl@europa:~/src/perl/testing$ cat -n contextual.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Contextual::Return;
4
5 sub sensible {
6 return STR { "one" }
7 NUM { 1 }
8 LIST { 1,2,3 }
9 HASHREF { {name => 'foo', value => 99} }
10 ;
11 }
```

```

12
13 print "Result = ".sensible()."\n";
14 print "Result = ".(0+sensible())."\n";
15 print "Result = ",sensible(),"\n";
16 print "Result = (name = ",sensible()->{name}," , value = ", sensible()->{value},")\n";

```

*Cuando se ejecuta, este programa produce la siguiente salida:*

```

pl@europa:~/src/perl/testing$./contextual.pl
Result = one
Result = 1
Result = 123
Result = (name = foo, value = 99)

```

*Las relaciones de definición-uso de la función **sensible** ¿Son un caso de binding estático o de binding dinámico? ¿Cuál es el ámbito de la declaración de **sensible** en las líneas 5-11?*

## Visibilidad

Como ya sabemos, es falso que en el ámbito de una declaración que define a **x** dicha declaración se aplique a todas las ocurrencias de **x** en su ámbito. En un ámbito estático, una declaración local a la anterior puede *ocultar la visibilidad* de la declaración anterior de **x**.

**Definición 12.1.8.** *Las reglas de visibilidad de un lenguaje especifican como se relacionan los nombres con las declaraciones que se les aplican.*

El concepto de nombre depende del lenguaje. Algunos lenguajes permiten que en un cierto ámbito haya mas de una definición asociada con un identificador. Un mecanismo que puede ser usado para determinar univocamente que definición se aplica a un determinado uso de un nombre es que el uso del nombre vaya acompañado de un sigil. (podría reinterpretarse que en realidad el nombre de una variable en Perl incluye el sigil).

Así, en Perl tenemos que es legal tener diferentes variables con nombre **x**: **\$x**, **@x**, **%x**, **&x**, **\*x**, etc. ya que van prefijadas por diferentes *sigils* **\$**, **@**, etc. El sigil que prefija **x** determina que definición se aplica al uso de **x** (La palabra sigil hace referencia a 'sellos mágicos' - combinaciones de símbolos y figuras geométricas - que son usados en algunas invocaciones con el propósito de producir un sortilegio).

En algunos casos existen mecanismos para la extensión de los nombres. En Perl es posible acceder a una variable de paquete ocultada por una léxica usando su nombre completo.

La asignación de una declaración a ciertos usos de un identificador puede requerir de otras fases de análisis semántico, como el análisis de tipos. El uso de los nombres de campo de un registro en Pascal y en C constituye un ejemplo:

```

1 type
2 a = ^b;
3 b = record
4 a: Integer;
5 b: Char;
6 c: a
7 end;
8 var
9 pointertob: a;
10 c : Integer;
11
12 ...
13 new(pointertob);
14
15 pointertob^.c := nil;
16 c := 4;
17 ...

```

El uso de `c` en la línea 15 es posible porque el tipo de la expresión `pointertob^` es un registro. La definición que se aplica al uso de `c` en la línea 16 es la de la línea 10.

También es posible hacer visible un nombre escondido - sin necesidad de extender el identificador - mediante alguna directiva que lo haga visible: Un ejemplo es la declaración `with the Pascal`:

```
new(pointertob);

with pointertob^ do
begin
 a := 10;
 b := 'A';
 c := nil
end;
...
```

### Declaraciones y Definiciones

En algunos lenguajes se distingue entre declaraciones que sólo proporcionan información sobre el elemento y lo hacen *visible* - pero no asignan memoria o producen código para la implementación del mismo - y otras que si producen dicho código. A las primeras se las suele llamar declaraciones y a las segundas definiciones. En tales lenguajes se considera un error que dos declaraciones de un mismo elemento difieran.

Por ejemplo, en C una variable o función sólo es definida una vez, pero puede ser declarada varias veces. El calificativo `extern` es usado en C para indicar que una declaración provee visibilidad pero no conlleva definición (creación):

```
extern char stack[10];
extern int stkptr;
```

Estas declaraciones le dicen al compilador C que las definiciones de los nombres `stack` y `stackptr` se encuentran en otro fichero. Si la palabra `extern` fuera omitida el compilador asignaría memoria para las mismas.

Otro ejemplo en el que una directiva hace visible una definición escondida es el uso de la declaración `our` de Perl cuando un paquete está repartido entre varios ficheros que usan repetitivamente `strict`:

**Ejercicio 12.1.5.** *Considere el siguiente programa:*

```
pl@europa:~/src/perl/perltesting$ cat -n useA.pl
1 #!/usr/bin/perl
2 package A;
3 use warnings;
4 use strict;
5
6 use A;
7
8 #our $x;
9 print "$x\n";
```

*La variable `$x` esta declarada en el fichero `A.pm`:*

```
pl@europa:~/src/perl/perltesting$ cat -n A.pm
1 package A;
2 use warnings;
3 use strict;
4
```

```

5 our $x = 1;
6
7 1;

```

*Sin embargo la compilación de useA.pl produce errores, pues \$x no es visible:*

```

pl@europa:~/src/perl/perltesting$ perl -c useA.pl
Variable "$x" is not imported at useA.pl line 9.
Global symbol "$x" requires explicit package name at useA.pl line 9.
useA.pl had compilation errors.

```

*El mensaje se arregla descomentando la declaración de \$x en la línea 8 de useA.pl:*

```

pl@europa:~/src/perl/perltesting$ perl -ce 'sed -e 's/#our/our/' useA.pl'
-e syntax OK

```

*La declaración de la línea 8 hace visible la variable \$x en el fichero useA.pl.*

*¿Cual es entonces el ámbito de la declaración de \$x en la línea 5 de A.pm? ¿Es todo el paquete? ¿O sólo el segmento del paquete que está en el fichero A.pm? (se supone que trabajamos con strict activado).*

## Inferencia, Declaraciones Implícitas y Ámbito

**Ejercicio 12.1.6.** *Tanto en los lenguajes estáticos como en los dinámicos se suele requerir que exista una declaración del objeto usado que determine las propiedades del mismo.*

*Aunque en los lenguajes dinámicos la creación/definición del elemento asociado con un nombre puede postergarse hasta el tiempo de ejecución, la generación de código para la sentencia de uso suele requerir un conocimiento (aunque sea mínimo) del objeto que esta siendo usado. En algunos casos, es necesario inferir la declaración a partir del uso, de manera que la declaración asociada con un uso es construida a partir del propio uso.*

*Los lenguajes típicamente estáticos fuertemente tipeados suelen requerir que para todo uso exista una declaración explícita y completa del nombre y de las operaciones que son válidas sobre el mismo al finalizar la fase de compilación. Sin embargo, el código para la creación/definición de algunos objetos puede ser postergado a la fase de enlace. De hecho, la resolución de ciertos enlaces pueden ocurrir durante la fase de ejecución (énlace dinámico).*

*El siguiente ejemplo hace uso de un typoglob selectivo en la línea 8 para definir la función ONE:*

```

pl@europa:~/src/perl/testing$ cat -n glob.pl
1 use warnings;
2 use strict;
3
4 sub one {
5 "1\n"
6 }
7
8 *ONE = \&one;
9 print ONE();

```

*Al ejecutar este programa se produce la salida;*

```

pl@europa:~/src/perl/testing$ perl glob.pl
1

```

*¿El uso de ONE en la línea 9 es un ejemplo de binding estático o dinámico? ¿Cuál es el ámbito de la declaración de ONE?*

*Responda estas mismas preguntas para esta otra variante del ejemplo anterior:*

```
pl@nereida:~/src/perl/perltesting$ cat -n coderef.pl
1 use warnings;
2 use strict;
3
4 *ONE = sub { "1\n" };
5 print ONE();
```

En los ejemplos anteriores el propio uso del nombre `ONE` actúa como una declaración: Perl deduce de la presencia de paréntesis después de `ONE` que `ONE` es el nombre de una función. Esta información es suficiente para generar el código necesario. Podría decirse que la forma del uso declara al ente `ONE` y que la línea de uso conlleva una declaración implícita. Sin embargo, la creación/definición completa de `ONE` es postergada hasta la fase de ejecución.

**Ejercicio 12.1.7.** *La conducta del compilador de Perl cambia si se sustituye el programa anterior por este otro:*

```
pl@europa:~/src/perl/testing$ cat -n globwarn.pl
1 use warnings;
2 use strict;
3
4 sub one {
5 "1\n"
6 }
7
8 *ONE = \&one;
9 my $x = ONE;
10 print $x;
```

*Al compilar se obtiene un error:*

```
pl@europa:~/src/perl/testing$ perl -c globwarn.pl
Bareword "ONE" not allowed while "strict subs" in use at globwarn.pl line 9.
globwarn.pl had compilation errors.
```

*¿Sabría explicar la causa de este cambio de conducta?*

**Ejercicio 12.1.8.** *El error que se observa en el ejercicio anterior desaparece cuando se modifica el código como sigue:*

```
lusasoft@LusaSoft:~/src/perl/perltesting$ cat -n globheader.pl
1 #!/usr/bin/perl
2 use warnings;
3 use strict;
4
5 sub ONE;
6
7 sub one {
8 "1\n"
9 }
10
11 *ONE = \&one;
12 my $x = ONE;
13 print $x;
```

*¿Cual es el significado de la línea 5?*

En el caso del lenguaje Simple C introducido en la práctica 12.3 hay una única declaración que se aplica a cada ocurrencia correcta de un nombre en el ámbito de dicha declaración.

Esto no tiene porque ser siempre así: en ciertos lenguajes una redeclaración de un cierto nombre *x* puede que sólo oculte a otra declaración previa de *x* si las dos declaraciones asignan el mismo tipo a *x*. Esta idea suele conocerse como *sobrecarga de identificadores*. De todos modos, sigue siendo cierto que para que el programa sea considerado correcto es necesario que sea posible *inferir* para cada ocurrencia de un identificador que única definición se aplica. Así una llamada a una cierta función `min(x,y)` llamaría a diferentes funciones `min` según fueran los tipos de *x* e *y*. Para resolver este caso es necesario combinar las fases de análisis de ámbito y de análisis de tipos.

Algunos lenguajes - especialmente los lenguajes funcionales - logran eliminar la mayoría de las declaraciones. Disponen de un mecanismo de inferencia que les permite - en tiempo de compilación - deducir del uso la definición y propiedades del nombre.

Véase como ejemplo de inferencia la siguiente sesión en OCaml

```
pl@nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples$ ocaml
Objective Caml version 3.09.2
```

```
let minimo = fun i j -> if i<j then i else j;;
val minimo : 'a -> 'a -> 'a = <fun>
minimo 2 3;;
- : int = 2
minimo 4.9 5.3;;
- : float = 4.9
minimo "hola" "mundo";;
- : string = "hola"
```

El compilador OCaml infiere el tipo de las expresiones. Así el tipo asociado con la función `minimo` es

```
'a -> 'a -> 'a
```

que es una *expresión de tipo* que contiene *variables de tipo*. El operador `->` es asociativo a derechas y así la expresión debe ser leída como `'a -> ('a -> 'a)`. Básicamente dice:

El tipo de la expresión es una función que toma un argumento de tipo `'a` (donde `'a` es una variable tipo que será instanciada en el momento del uso de la función) y devuelve una función que toma elementos de tipo `'a` y retorna elementos de tipo `'a`.

## Ámbito Dinámico

En el *ámbito dinámico*, cada nombre para el que se usa ámbito dinámico tiene asociada una pila de bindings. Cuando se crea un nuevo ámbito dinámico se empuja en la pila el viejo valor (que podría no estar definido). Cuando se sale del ámbito se saca de la pila el antiguo valor. La evaluación de *x* retorna siempre el valor en el top de la pila.

La sentencia `local` de Perl provee de ámbito dinámico a las variables de paquete. Una aproximación a lo que ocurre cuando se ejecuta `local` es:

| DECLARACIÓN DE <code>local</code>                                  | SIGNIFICADO                                                                                                                |
|--------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|
| <pre>{   local(\$SomeVar);   \$SomeVar = 'My Value';   ... }</pre> | <pre>{   my \$TempCopy = \$SomeVar;   \$SomeVar = undef;   \$SomeVar = 'My Value';   ...   \$SomeVar = \$TempCopy; }</pre> |

La diferencia entre ámbito dinámico y estático debería quedar mas clara observando la conducta del siguiente código

```
lhp@nereida:~/Lperl/src$ cat -n local.pl
1 #!/usr/bin/perl -w
2 use strict;
3
4 our $x;
5
6 sub pr { print "$x\n"; }
7 sub titi { my $x = "titi"; pr(); }
8 sub toto { local $x = "toto"; &pr(); &titi(); }
9
10 $x = "global";
11 &pr();
12 &toto();
13 &titi();
```

Cuando se ejecuta, se obtiene la siguiente salida:

```
> local.pl
global
toto
toto
global
```

**Ejercicio 12.1.9.** *¿Es local una declaración o una sentencia? ¿Que declaraciones se aplican a los diferentes usos de \$x en las líneas 6, 7, 8 y 10?*

## 12.2. Descripción Eyapp del Lenguaje SimpleC

El proceso de identificar los nombres conlleva establecer enlaces entre las ocurrencias y sus declaraciones o bien - en caso de error - determinar que dicho enlace no existe. El resultado de este proceso de identificación (o análisis de ámbito y visibilidad) será utilizado durante las fases posteriores.

En este capítulo usaremos Parse::Eyapp para desarrollar las primeras fases - análisis léxico, análisis sintáctico y análisis de ámbito - de un compilador para un subconjunto de C al que denominaremos Simple C :

### El Cuerpo

```
%%
program:
 definition<%name PROGRAM +>.program
 { $program }
;

definition:
 $funcDef
 { $funcDef }
| %name FUNCTION
 $basictype $funcDef
 { $funcDef }
| declaration
 {}
```

```

;

basicType:
 %name INT
 INT
 | %name CHAR
 CHAR
;

funcDef:
 $ID '(' $params ')' $block
 {
 flat_statements($block);
 $block->{parameters} = [];
 $block->{function_name} = $ID;
 $block->type("FUNCTION");
 return $block;
 }
;

params:
 (basicType ID arraySpec) <%name PARAMS * ','>
 { $_[1] }
;

block:
 '{'.bracket
 declaration<%name DECLARATIONS *>.decs statement<%name STATEMENTS *>.sts '}'
 {
 flat_statements($sts);
 $sts->type("BLOCK") if $decs->children;
 return $sts;
 }
;

declaration:
 %name DECLARATION
 $basicType $declList ';'
;

declList:
 (ID arraySpec) <%name VARLIST + ','> { $_[1] }
;

arraySpec:
 ('[' INUM ']') * { $_[1]->type("ARRAYSPEC"); $_[1] }
;

statement:
 expression ';' { $_[1] }
 | ';'
 | %name BREAK
 $BREAK ';'

```



```

| %name CONTINUE
 $CONTINUE ';'
| %name EMPTYRETURN
 RETURN ';'
| %name RETURN
 RETURN expression ';'
| block { $_[1] }
| %name IF
 ifPrefix statement %prec '+'
| %name IFELSE
 ifPrefix statement 'ELSE' statement
| %name WHILE
 $loopPrefix statement
;

ifPrefix:
 IF '(' expression ')' { $_[3] }
;

loopPrefix:
 $WHILE '(' expression ')' { $_[3] }
;

expression:
 binary <+ ',>'
 {
 return $_[1]->child(0) if ($_[1]->children() == 1);
 return $_[1];
 }
;

Variable:
 %name VAR
 ID
| %name VARARRAY
 $ID ('[' binary ']') <%name INDEXSPEC +>
;

Primary:
 %name INUM
 INUM
| %name CHARCONSTANT
 CHARCONSTANT
| $Variable
 {
 return $Variable
 }
| '(' expression ')' { $_[2] }
| $function_call
 {
 return $function_call # bypass
 }
;

```

```
function_call:
 %name
 FUNCTIONCALL
 ID '(' binary <%name ARGLIST * ','> ')';
;
```

```
Unary:
 '++' Variable
 | '--' Variable
 | Primary { $_[1] };
;
```

```
binary:
 Unary { $_[1] }
 | %name PLUS
 binary '+' binary
 | %name MINUS
 binary '-' binary
 | %name TIMES
 binary '*' binary
 | %name DIV
 binary '/' binary
 | %name MOD
 binary '%' binary
 | %name LT
 binary '<' binary
 | %name GT
 binary '>' binary
 | %name GE
 binary '>=' binary
 | %name LE
 binary '<=' binary
 | %name EQ
 binary '==' binary
 | %name NE
 binary '!=' binary
 | %name AND
 binary '&' binary
 | %name EXP
 binary '**' binary
 | %name OR
 binary '|' binary
 | %name ASSIGN
 $Variable '=' binary
 | %name PLUSASSIGN
 $Variable '+=' binary
 | %name MINUSASSIGN
 $Variable '-=' binary
 | %name TIMESASSIGN
 $Variable '*=' binary
 | %name DIVASSIGN
 $Variable '/=' binary
```

```

 | %name MODASSIGN
 $Variable '%=' binary
;

```

## La Cabeza

```

/*
File: Simple/Syntax.eypp
Full Type checking
To build it, Do make or:
 eyapp -m Simple::Syntax Syntax.eypp;
*/
%{
use strict;
use Carp;
use warnings;
use Data::Dumper;
use List::MoreUtils qw(firstval);
our $VERSION = "0.4";

my $debug = 1;
my %reserved = (
 int => "INT",
 char => "CHAR",
 if => "IF",
 else => "ELSE",
 break => "BREAK",
 continue => "CONTINUE",
 return => "RETURN",
 while => "WHILE"
);

my %lexeme = (
 '=' => "ASSIGN",
 '+' => "PLUS",
 '-' => "MINUS",
 '*' => "TIMES",
 '/' => "DIV",
 '%' => "MOD",
 '|' => "OR",
 '&' => "AND",
 '{' => "LEFTKEY",
 '}' => "RIGHTKEY",
 ',' => "COMMA",
 ';' => "SEMICOLON",
 '(' => "LEFTPARENTHESIS",
 ')' => "RIGHTPARENTHESIS",
 '[' => "LEFTBRACKET",
 ']' => "RIGHTBRACKET",
 '==' => "EQ",
 '+=' => "PLUSASSIGN",
 '-=' => "MINUSASSIGN",
 '*=' => "TIMESASSIGN",
 '/=' => "DIVASSIGN",

```

```

'%= ' => "MODASSIGN",
'!= ' => "NE",
'< ' => "LT",
'> ' => "GT",
'<=' => "LE",
'>=' => "GE",
'++ ' => "INC",
'-- ' => "DEC",
'** ' => "EXP"
);

my ($tokenbegin, $tokenend) = (1, 1);

sub flat_statements {
 my $block = shift;

 my $i = 0;
 for ($block->children) {
 if ($_->type eq "STATEMENTS") {
 splice @{$block->{children}}, $i, 1, $_->children;
 }
 $i++;
 }
}
%}

%syntactic token RETURN BREAK CONTINUE

%right '= ' '+=' '-=' '*=' '/=' '%='
%left '|'
%left '&'
%left '==' '!='
%left '<' '>' '>=' '<='
%left '+' '-'
%left '*' '/' '%'
%right '**'
%right '++' '--'
%right 'ELSE'

%tree

La Cola

%%

sub _Error {
 my($token)=$_[0]->YYCurval;
 my($what)= $token ? "input: '$token->[0] ' in line $token->[1]" : "end of input";
 my @expected = $_[0]->YYExpect();
 my $expected = @expected? "Expected one of these tokens: '@expected':"";

 croak "Syntax error near $what. $expected\n";
}

```

```

sub _Lexer {
 my($parser)=shift;

 my $token;
 for ($parser->YYData->{INPUT}) {
 return('','undef) if !defined($_) or $_ eq '';

 #Skip blanks
 s{\A
 ((?:
 \s+ # any white space char
 | /*.*?*/ # C like comments
)+
)
 }
 {}xs
 and do {
 my($blanks)=$1;

 #Maybe At EOF
 return('','undef) if $_ eq '';
 $tokenend += $blanks =~ tr/\n//;
 };

 $tokenbegin = $tokenend;

 s/^('.')//
 and return('CHARCONSTANT', [$1, $tokenbegin]);

 s/^(([0-9]+(?:\.[0-9]+)?))//
 and return('INUM',[$1, $tokenbegin]);

 s/^(([A-Za-z][A-Za-z0-9_]*)//
 and do {
 my $word = $1;
 my $r;
 return ($r, [$r, $tokenbegin]) if defined($r = $reserved{$word});
 return('ID',[$word, $tokenbegin]);
 };

 m/^(\\S\\S)/ and defined($token = $1) and exists($lexeme{$token})
 and do {
 s/\\.//;
 return ($token, [$token, $tokenbegin]);
 }; # do

 m/^(\\S)/ and defined($token = $1) and exists($lexeme{$token})
 and do {
 s/\\.//;
 return ($token, [$token, $tokenbegin]);
 }; # do

 die "Unexpected character at $tokenbegin\n";
 }
}

```

```

 } # for
}

sub compile {
 my($self)=shift;

 my ($t);

 $self->YYData->{INPUT} = $_[0];

 $t = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
 #yydebug => 0x1F
);

 return $t;
}

sub TERMINAL::value {
 return $_[0]->{attr}[0];
}

line Support

sub TERMINAL::line {
 return $_[0]->{attr}[1];
}

sub VAR::line {
 my $self = shift;

 return $self->child(0)->{attr}[1];
}

sub PLUS::line {
 $_[0]->{lines}[0]
}

{
no warnings;
*TIMES::line = *DIV::line = *MINUS::line = *ASSIGN::line
=*GT::line
=*IF::line
=*RETURN::line
= \&PLUS::line;

*VARARRAY::line = *FUNCTIONCALL::line
=\&VAR::line;
}

Debugging and Display
sub show_trees {
 my ($t) = shift;
 my $debug = shift;

```

```

$Data::Dumper::Indent = 1;
print Dumper $t if $debug > 3;
local $Parse::Eyapp::Node::INDENT = $debug;
print $t->str."\n";
}

sub TERMINAL::info {
 my $a = join ':', @{$_[0]->{attr}};
 return $a
}

sub TERMINAL::save_attributes {
 # $_[0] is a syntactic terminal
 # $_[1] is the father.
 push @{$_[1]->{lines}}, $_[0]->[1]; # save the line!
}

sub WHILE::line {
 return $_[0]->{line}
}

```

### Ejemplo de Árbol Construido

pl@nereida:~/Lbook/code/Simple-Syntax/script\$ usesyntax.pl bugmatch01.c 2

test (int n)

```

{
 while (1) {
 if (1>0) {
 a = 2;
 break;
 }
 else if (2> 0){
 b = 3;
 continue;
 }
 }
}

```

```

PROGRAM(
 FUNCTION(
 WHILE(
 INUM(
 TERMINAL[1:3]
),
 STATEMENTS(
 IFELSE(
 GT(
 INUM(
 TERMINAL[1:4]
),
 INUM(
 TERMINAL[0:4]
)
)
)
)
)
)
)

```

```

) # GT,
STATEMENTS(
 ASSIGN(
 VAR(
 TERMINAL[a:5]
),
 INUM(
 TERMINAL[2:5]
)
) # ASSIGN,
 BREAK
) # STATEMENTS,
IF(
 GT(
 INUM(
 TERMINAL[2:8]
),
 INUM(
 TERMINAL[0:8]
)
) # GT,
 STATEMENTS(
 ASSIGN(
 VAR(
 TERMINAL[b:9]
),
 INUM(
 TERMINAL[3:9]
)
) # ASSIGN,
 CONTINUE
) # STATEMENTS
) # IF
) # IFELSE
) # STATEMENTS
) # WHILE
) # FUNCTION
) # PROGRAM

```

### Ejemplo de Árbol con Aplanamiento de STATEMENTS

```

pl@nereida:~/Lbook/code/Simple-Syntax/script$ usesyntax.pl prueba26.c 2
int a[20],b,e[10];

```

```

g() {}

```

```

int f(char c) {
char d;
c = 'X';
e[d][b] = 'A'+c;
{
 int d;
 d = c * 2;
}
}

```



```

{
 d = a + b;
 {
 c = a + 1;
 }
}
c = d * 2;
return c;
}

```

```

PROGRAM(
 FUNCTION,
 FUNCTION(
 ASSIGN(
 VAR(
 TERMINAL[c:7]
),
 CHARCONSTANT(
 TERMINAL['X':7]
)
) # ASSIGN,
 ASSIGN(
 VARARRAY(
 TERMINAL[e:8],
 INDEXSPEC(
 VAR(
 TERMINAL[d:8]
),
 VAR(
 TERMINAL[b:8]
)
) # INDEXSPEC
) # VARARRAY,
 PLUS(
 CHARCONSTANT(
 TERMINAL['A':8]
),
 VAR(
 TERMINAL[c:8]
)
) # PLUS
) # ASSIGN,
 BLOCK(
 ASSIGN(
 VAR(
 TERMINAL[d:11]
),
 TIMES(
 VAR(
 TERMINAL[c:11]
),
 INUM(

```

```

 TERMINAL[2:11]
)
) # TIMES
) # ASSIGN
) # BLOCK,
ASSIGN(
 VAR(
 TERMINAL[d:14]
),
 PLUS(
 VAR(
 TERMINAL[a:14]
),
 VAR(
 TERMINAL[b:14]
)
) # PLUS
) # ASSIGN,
ASSIGN(
 VAR(
 TERMINAL[c:16]
),
 PLUS(
 VAR(
 TERMINAL[a:16]
),
 INUM(
 TERMINAL[1:16]
)
) # PLUS
) # ASSIGN,
ASSIGN(
 VAR(
 TERMINAL[c:19]
),
 TIMES(
 VAR(
 TERMINAL[d:19]
),
 INUM(
 TERMINAL[2:19]
)
) # TIMES
) # ASSIGN,
RETURN(
 VAR(
 TERMINAL[c:20]
)
) # RETURN
) # FUNCTION
) # PROGRAM

```

## 12.3. Práctica: Construcción del AST para el Lenguaje Simple C

Reproduzca el analizador sintáctico para el lenguaje SimpleC introducido en la sección 12.2.

## 12.4. Práctica: Análisis de Ámbito del Lenguaje Simple C

Haga una primera parte del análisis de ámbito del lenguaje Simple C presentado en la práctica 12.3.

En esta primera parte se construyen las declaraciones y las tablas de símbolos de cada bloque. Sin embargo no queda establecida la asignación desde una instancia de un objeto a su declaración.

**La tabla de Símbolos** Comenzaremos describiendo las estructuras de datos que van a conformar la tabla de símbolos de nuestro compilador:

1. Dote de los siguientes atributos a los nodos de tipo "bloque", esto es a aquellos nodos que definen ámbito. Por ejemplo, los nodos asociados con la variable sintáctica `block` son de tipo "bloque":

```
block:
 %name BLOCK
 '{' declaration %name DECLARATIONS * statement %name STATEMENTS * '}'
```

- a) Atributo `symboltable`: referencia a la tabla de símbolos asociada con el bloque La tabla de símbolos asociada al bloque es un simple hash.
  - b) Atributo `fatherblock`: referencia al nodo `BLOCK` que inmediatamente anida a este. Los nodos de tipo bloque quedan enlazados según el árbol de anidamiento de bloques
  - c) Los bloques con declaraciones vacías pueden ser simplificados en el árbol de bloques
2. Los nodos `FUNCTION` asociados con las funciones son nodos de tipo bloque y serán tratados de manera similar a los nodos `BLOCK`, esto es, tendrán su tabla de símbolos asociada en la cual se guardarán los parámetros de la función. Este bloque es el bloque padre del bloque formado por el cuerpo de la función. Posteriormente se pueden fusionar estos dos bloques siempre que se conserve la información necesaria sobre los parámetros.

```
funcDef:
 %name FUNCTION
 ID '(' param <%name PARAMS * ',' '>' '>' ')'
 block
```

3. Los identificadores de funciones van en la tabla de símbolos global, asociada con el nodo `PROGRAM`:

```
pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^program:/,/;/p' SymbolTables
1 program:
2 { reset_file_scope_vars(); }
3 definition<%name PROGRAM +>.program
4 {
5 $program->{symboltable} = { %st }; # creates a copy of the s.t.
6 for (keys %type) {
7 $type{$_} = Parse::Eyapp::Node->hnew($_);
8 }
9 $program->{depth} = 0;
10 $program->{line} = 1;
11 $program->{types} = { %type };
12 $program->{lines} = $tokenend;
```

```

13
14 reset_file_scope_vars();
15 $program;
16 }
17 ;

```

Observe que en C una función puede ser usada antes de que aparezca su definición.

4. Las instancias de los objetos y en particular los nodos **VAR** deben tener atributos que determinen que declaración se les aplica y en que ámbito viven. Sin embargo estos atributos no serán trabajados en esta práctica:

- a) Atributo **entry**: Su entrada en la tabla de símbolos a la que pertenece
- b) Atributo **scope**: Una referencia al nodo **BLOCK** en el que fué declarada la variable

5. La tabla de símbolos es un árbol de tablas. Cada tabla esta relacionada con un bloque. Cada tabla tiene una entrada para cada identificador que fué declarado en su bloque. En dicha entrada figuran los atributos del identificador: entre estos últimos la línea en la que fue declarado y su tipo.

**¿Que es una declaración?** Para facilitar el análisis de ámbito y la comprobación de tipos debemos modificar la fase de construcción del AST para producir declaraciones que permitan comprobar con facilidad la *equivalencia de tipos*

La equivalencia de tipos se realiza habitualmente mediante *expresiones de tipo*. En nuestro caso vamos a elegir una representación dual mediante cadenas y árboles de las expresiones de tipo. Las cadenas serán términos árbol. Como primera aproximación, entenderemos que dos tipos son equivalentes si las cadenas que representan sus expresiones de tipo son iguales.

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '1,/%}/p' SymbolTables.eyyp | c
1 /*
2 File: SymbolTables.eyyp
3 Scope Analysis: Only symbol table construction
4 */
5 %{
6 use strict;
7 use Data::Dumper;
8 use List::MoreUtils qw(firstval lastval);
9 use Simple::Trans;
10
11 my %reserved = (
12 => "....."
13);
14
15 my %lexeme = (
16 '..' => ".."
17);
18
19 our ($blocks);
20
21 sub is_duplicated {
22
23 }
24
25 sub build_type {

```

```

..
82 }
83
84 my ($tokenbegin, $tokenend);
85 my %type;
86 my $depth;
87
88 my %st; # Global symbol table
89
90 sub build_function_scope {
..
114 }
115
116 sub reset_file_scope_vars {
117 %st = (); # reset symbol table
118 ($tokenbegin, $tokenend) = (1, 1);
119 %type = (INT => 1,
120 CHAR => 1,
121 VOID => 1,
122);
123 $depth = 0;
124 }
125
126 %}

```

**El código asociado con declaration** Veamos como modificamos la construcción del AST durante las declaraciones:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^declaration:/,/^;/p' Symbol
1 declaration:
2 %name DECLARATION
3 $basicstype $declList ','
4 {
5 my %st; # Symbol table local to this declaration
6 my $bt = $basicstype->type;
7 my @decs = $declList->children();
8 while (my ($id, $arrspec) = splice(@decs, 0, 2)) {
9 my $name = $id->{attr}[0];
10 my $type = build_type($bt, $arrspec);
11 $type{$type} = 1; # has too much $type for me!
12
13 # control duplicated declarations
14 die "Duplicated declaration of $name at line $id->{attr}[1]\n" if exists($st{$name});
15 $st{$name}->{type} = $type;
16 $st{$name}->{line} = $id->{attr}[1];
17 }
18 return \%st;
19 }
20 ;

```

El código de la función `build_type` es como sigue:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^sub build_type/,82p' Symbol
1 sub build_type {

```

```

2 my $bt = shift;
3 my @arrayspec = shift()->children();
4
5 my $type = '';
6 for my $s (@arrayspec) {
7 $type .= "A_$s->{attr}[0](";
8 }
9 if ($type) {
10 $type = "$type$bt"."(")x@arrayspec);
11 }
12 else {
13 $type = $bt;
14 }
15 return $type;
16 }

```

**Tratamiento de los Bloques** La variable sintáctica `bloque` genera el lenguaje de las definiciones seguidas de sentencias:

```
'{ { datadefinition } { statement } }'
```

El lenguaje de los bloques es modificado en consonancia:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^block:\/,\/^\/p' SymbolTables
1 block:
2 '{'.bracket declaration<%name DECLARATIONS *>.decs statement<%name STATEMENTS *>.sts
3 {
4 my %st;
5
6 for my $lst ($decs->children) {
7
8 # control duplicated declarations
9 my $message;
10 die $message if $message = is_duplicated(\%st, $lst);
11
12 %st = (%st, %$lst);
13 }
14 $sts->{symboltable} = \%st;
15 $sts->{line} = $bracket->[1];
16 $sts->type("BLOCK") if %st;
17 return $sts;
18 }
19
20 ;

```

El código de la función `is_duplicated` es como sigue:

```

sub is_duplicated {
 my ($st1, $st2) = @_;

 my $id;

 defined($id=firstval{exists $st1->{$_}} keys %$st2)
 and return "Error. Variable $id at line $st2->{$id}->{line} declared twice.\n";
 return 0;
}

```

**Tratamiento de las Funciones** El código para las funciones es similar. Unimos la tabla de símbolos de los parámetros a la del bloque para posteriormente reconvertir el bloque en un nodo **FUNCTION**:

```
pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^funcDef:/,/;/p' SymbolTable
1 funcDef:
2 $ID '(' $params ')'
3 $block
4 {
5 my $st = $block->{symboltable};
6 my @decs = $params->children();
7 $block->{parameters} = [];
8 while (my ($bt, $id, $arrspec) = splice(@decs, 0, 3)) {
9 my $bt = ref($bt); # The string 'INT', 'CHAR', etc.
10 my $name = $id->{attr}[0];
11 my $type = build_type($bt, $arrspec);
12 $type{$type} = 1; # has too much $type for me!
13
14 # control duplicated declarations
15 #die "Duplicated declaration of $name at line $id->{attr}[1]\n" if exists($st->
16 die "Duplicated declaration of $name at line $st->{$name}{line}\n" if exists($
17 $st->{$name}->{type} = $type;
18 $st->{$name}->{param} = 1;
19 $st->{$name}->{line} = $id->{attr}[1];
20 push @{$block->{parameters}}, $name;
21 }
22 $block->{function_name} = $ID;
23 $block->type("FUNCTION");
24 return $block;
25 }
26 ;
```

En la línea 8 se procesan las declaraciones de parámetros. La lista de parámetros estaba definida por:

```
params:
 (basictype ID arraySpec)<%name PARAMS * ', '>
 { $_[1] }
```

En una segunda fase, a la altura de **definition** se construye la entrada para la función en la tabla de símbolos global:

```
pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^definition:/,/;/p' SymbolTable
1 definition:
2 $funcDef
3 {
4 build_function_scope($funcDef, 'INT');
5 }
6 | %name FUNCTION
7 $basictype $funcDef
8 {
9 build_function_scope($funcDef, $basictype->type);
10 }
11 | declaration
12 {
13 #control duplicated declarations
```

```

14 my $message;
15 die $message if $message = is_duplicated(\%st, $_[1]);
16 %st = (%st, %{$_[1]}); # improve this code
17 return undef; # will not be inserted in the AST
18 }
19 ;

```

El código de la función `build_function_scope` es como sigue:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/sub build_func/,114p' Symbol
1 sub build_function_scope {
2 my ($funcDef, $returntype) = @_;
3
4 my $function_name = $funcDef->{function_name}[0];
5 my @parameters = @{$funcDef->{parameters}};
6 my $lst = $funcDef->{symboltable};
7 my $numargs = scalar(@parameters);
8
9 #compute type
10 my $partype = "";
11 my @types = map { $lst->{$_}{type} } @parameters;
12 $partype .= join ",", @types if @types;
13 my $type = "F(X_$numargs($partype),$returntype)";
14
15 #insert it in the hash of types
16 $type{$type} = 1;
17
18 #insert it in the global symbol table
19 die "Duplicated declaration of $function_name at line $funcDef-->{attr}[1]\n"
20 if exists($st{$function_name});
21 $st{$function_name}->{type} = $type;
22 $st{$function_name}->{line} = $funcDef->{function_name}[1]; # redundant
23
24 return $funcDef;
25 }

```

**El Tratamiento del Programa Principal** Sigue el código (incompleto) asociado con la variable sintáctica de arranque:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$ sed -ne '/^program:/,/~/p' SymbolTabl
1 program:
2 { reset_file_scope_vars(); }
3 definition<%name PROGRAM +>.program
4 {
5 $program->{symboltable} = { %st }; # creates a copy of the s.t.
6 for (keys %type) {
7 $type{$_} = Parse::Eyapp::Node->hnew($_);
8 }
9 $program->{depth} = 0;
10 $program->{line} = 1;
11 $program->{types} = { %type };
12 $program->{lines} = $tokenend;
13
14 reset_file_scope_vars();

```



```

15 $program;
16 }
17 ;

```

Las árboles asociados con las expresiones de tipo son calculados en las líneas 6-8.

### Ejemplo de Salida: Control de declaraciones duplicadas

A estas alturas del análisis de ámbito podemos controlar la aparición de declaraciones duplicadas: para el programa de entrada:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/script$ usesymboltables.pl declaredtwice.c 2
int a, b[1][2];
char d, e[1][2];
char f(int a, char b[10]) {
 int c[1][2];
 char b[10], e[9]; /* b is declared twice */

 return b[0];
}

```

Duplicated declaration of b at line 5

Sin embargo no podemos conocer los objetos no declarados (variables, funciones) hasta que hayamos finalizado el análisis de ámbito.

### Un ejemplo de Volcado de un AST

Sigue un ejemplo de salida para un programa sin errores de ámbito:

```

pl@nereida:~/Lbook/code/Simple-SymbolTables/script$ usesymboltables.pl prueba10.c 2
int a,b;
f(int a, int b) {
 b = a+1;
}

int main() {
 a = 4;
 {
 int d;
 d = f(a);
 }
}

PROGRAM^{0}{(
 FUNCTION[f]^{1}{(
 ASSIGN(
 VAR[No declarado!](
 TERMINAL[b:3]
),
 PLUS(
 VAR[No declarado!](
 TERMINAL[a:3]
),
 INUM(
 TERMINAL[1:3]

```

```

)
) # PLUS
) # ASSIGN
) # FUNCTION,
FUNCTION[main]^{2}(
 ASSIGN(
 VAR[No declarado!](
 TERMINAL[a:7]
),
 INUM(
 TERMINAL[4:7]
)
) # ASSIGN,
BLOCK[8]^{3}(
 ASSIGN(
 VAR[No declarado!](
 TERMINAL[d:10]
),
 FUNCTIONCALL[No declarado!](
 TERMINAL[f:10],
 ARGLIST(
 VAR[No declarado!](
 TERMINAL[a:10]
)
) # ARGLIST
) # FUNCTIONCALL
) # ASSIGN
) # BLOCK
) # FUNCTION
) # PROGRAM

```

-----

0)

Types:

```

$VAR1 = {
 'CHAR' => bless({
 'children' => []
 }, 'CHAR'),
 'VOID' => bless({
 'children' => []
 }, 'VOID'),
 'INT' => bless({
 'children' => []
 }, 'INT'),
 'F(X_0(),INT)' => bless({
 'children' => [
 bless({
 'children' => []
 }, 'X_0'),
 $VAR1->{'INT'}
]
 }, 'F'),
 'F(X_2(INT,INT),INT)' => bless({
 'children' => [

```

```

 bless({
 'children' => [
 $VAR1->{'INT'},
 $VAR1->{'INT'}
]
 }, 'X_2'),
 $VAR1->{'INT'}
]
}, 'F')
};
Symbol Table of the Main Program:
$VAR1 = {
 'a' => {
 'type' => 'INT',
 'line' => 1
 },
 'b' => {
 'type' => 'INT',
 'line' => 1
 },
 'f' => {
 'type' => 'F(X_2(INT,INT),INT)',
 'line' => 2
 },
 'main' => {
 'type' => 'F(X_0(),INT)',
 'line' => 6
 }
};

```

---

```

1)
Symbol Table of f
$VAR1 = {
 'a' => {
 'type' => 'INT',
 'param' => 1,
 'line' => 2
 },
 'b' => {
 'type' => 'INT',
 'param' => 1,
 'line' => 2
 }
};

```

---

```

2)
Symbol Table of main
$VAR1 = {};

```

---

```

3)

```

Symbol Table of block at line 8

```
$VAR1 = {
 'd' => {
 'type' => 'INT',
 'line' => 9
 }
};
```

pl@nereida:~/Lbook/code/Simple-SymbolTables/script\$

## Métodos de soporte para el Volcado del Arbol

Se han utilizado las siguientes funciones de soporte para producir el volcado:

pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple\$ sed -ne '683,\$p' SymbolTables.eyp

```
Debugging and Display
```

```
sub show_trees {
 my ($t) = @_;
```

```

 print Dumper $t if $debug > 1;
 $Parse::Eyapp::Node::INDENT = 2;
 $Data::Dumper::Indent = 1;
 print $t->str."\n";
}
```

```

$Parse::Eyapp::Node::INDENT = 1;
sub TERMINAL::info {
 my @a = join ':', @{$_[0]->{attr}};
 return "@a"
}
```

```

sub PROGRAM::footnote {
 return "Types:\n"
 .Dumper($_[0]->{types}).
 "Symbol Table of the Main Program:\n"
 .Dumper($_[0]->{symboltable})
}
```

```

sub FUNCTION::info {
 return $_[0]->{function_name}[0]
}
```

```

sub FUNCTION::footnote {
 my $text = '';
 $text .= "Symbol Table of $_[0]->{function_name}[0]\n" if $_[0]->type eq 'FUNCTION';
 $text .= "Symbol Table of block at line $_[0]->{line}\n" if $_[0]->type eq 'BLOCK';
 $text .= Dumper($_[0]->{symboltable});
 return $text;
}
```

```

sub BLOCK::info {
 my $info = "$_[0]->{line}";
 return $info;
}
```

```

*BLOCK::footnote = \&FUNCTION::footnote;

sub VAR::info {
 return $_[0]->{type} if defined $_[0]->{type};
 return "No declarado!";
}

*FUNCTIONCALL::info = *VARARRAY::info = \&VAR::info;
pl@nereida:~/Lbook/code/Simple-SymbolTables/lib/Simple$

```

## 12.5. La Dificultad de Elaboración de las Pruebas

A la hora de hacer las pruebas necesitamos comprobar que dos árboles son iguales. El problema que aparece en el diseño de un compilador es un caso particular de la regla *defina su API antes de realizar las pruebas* y de la ausencia de herramientas adecuadas.

Como el diseño de un compilador se hace por fases y las fases a veces se superponen resulta que el resultado de una llamada a la misma subrutina cambia en las diversas etapas del desarrollo. Terminada la fase de análisis sintáctico el producto es un AST. cuando posteriormente añadimos la fase de análisis de ámbito obtenemos un árbol decorado. Si hicimos pruebas para la primera fase y hemos usado `is_deeply` las pruebas no funcionarían con la versión ampliada a menos que se cambie el árbol esperado. Ello es debido a que `is_deeply` requiere la igualdad estructural total de los dos árboles.

### Limitaciones de Data::Dumper

La descripción dada con `Data::Dumper` de una estructura de datos anidada como es el árbol es difícil de seguir, especialmente como en este ejemplo en el que hay enlaces que autoreferencian la estructura.

```

$VAR1 = bless({
 'types' => {
 'CHAR' => bless({ 'children' => [] }, 'CHAR'),
 'INT' => bless({ 'children' => [] }, 'INT'),
 'F(X_0(),CHAR)' => bless({ 'children' => [
 bless({ 'children' => [] }, 'X_0'), bless({ 'children' => [] }, 'CHAR')
]
 }, 'F'),
 'F(X_0(),INT)' => bless({ 'children' => [
 bless({ 'children' => [] }, 'X_0'), bless({ 'children' => [] }, 'INT')
]
 }, 'F')
 },
 'symboltable' => {
 'd0' => { 'type' => 'CHAR', 'line' => 1 },
 'f' => { 'type' => 'F(X_0(),CHAR)', 'line' => 2 },
 'g' => { 'type' => 'F(X_0(),INT)', 'line' => 13 },
 },
 'lines' => 19,
 'children' => [
 bless({ # FUNCTION f <-----x
 'parameters' => [], 'symboltable' => {}, 'fatherblock' => $VAR1,
 'function_name' => ['f', 2],
 'children' => [
 bless({ <-----|-----x

```

```

'line' => 3
'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0], -----x
'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0]['children'][0],
 'children' => [], 'line' => 4
 }, 'BLOCK')
],
}, 'BLOCK'),
bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0], -----x
 'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0]['children'][1],
 'children' => [], 'line' => 7
 }, 'BLOCK')
],
 'line' => 6
}, 'BLOCK'),
bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0],
 'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0]['children'][2],
 'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[0]['children'}[2]{
 'children' => [], 'line' => 10
 }, 'BLOCK')
],
 'line' => 10
 }, 'BLOCK')
],
 'line' => 9
}, 'BLOCK')
],
'line' => 2
}, 'FUNCTION'),
bless({ # FUNCTION g
 'parameters' => [], 'symboltable' => {}, 'fatherblock' => $VAR1,
 'function_name' => ['g', 13],
 'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[1],
 'children' => [], 'line' => 14
 }, 'BLOCK'),
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[1],
 'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[1]['children'}[1],
 'children' => [], 'line' => 16
 }, 'BLOCK')
],

```

```

],
 'line' => 15
 }, 'BLOCK'),
 bless({
 'symboltable' => {}, 'fatherblock' => $VAR1->{'children'}[1],
 'children' => [], 'line' => 18
 }, 'BLOCK')
],
'line' => 13
}, 'FUNCTION')
],
'line' => 1
}, 'PROGRAM');

```

La estructura es en realidad el árbol de análisis abstracto decorado para un programa SimpleC (véase el capítulo 12). Para ver el fuente que se describe y la estructura del AST consulte la tabla 12.1. Al decorar el árbol con la jerarquía de bloques (calculada en la sección 10.4, página ??) se producen numerosas referencias cruzadas que hacen difícil de leer la salida de `Data::Dumper`.

### La opción `Data::Dumper::Purity`

El problema no es sólo que es difícil seguir una estructura que se autoreferencia como la anterior. Una estructura recursiva como esta no puede ser evaluada como código Perl, ya que *Perl prohíbe que una variable pueda ser usada antes de que finalice su definición*. Por ejemplo, la declaración-inicialización:

```
my $x = $x -1;
```

es considerada incorrecta. Eso significa que *el texto producido por Data::Dumper de esta forma no puede ser insertado en un test de regresión*.

El módulo `Data::Dumper` dispone de la variable `Data::Dumper::Purity` la cual ayuda a subsanar esta limitación. Veamos la salida que se produce para el programa anterior cuando `Data::Dumper::Purity` esta activa:

```

$VAR1 = bless({
 'types' => { ... },
 'symboltable' => { ... },
 'lines' => 19,
 'children' => [
 bless({
 'function_name' => ['f', 2], 'line' => 2,
 'parameters' => [], 'symboltable' => {}, 'fatherblock' => {},
 'children' => [
 bless({
 'symboltable' => {}, 'fatherblock' => {}, 'line' => 3
 'children' => [
 bless({ 'symboltable' => {}, 'fatherblock' => {}, 'children' => [], 'line' => 4 }
],
 }, 'BLOCK'),
],
 }, 'BLOCK'),
],
}, 'BLOCK'),
bless({
 'symboltable' => {}, 'fatherblock' => {}, 'line' => 6
 'children' => [
 bless({
 'symboltable' => {},
 'fatherblock' => {},
 'children' => [],

```

```

 'line' => 7
 }, 'BLOCK')
],
}, 'BLOCK'),
bless({
 'line' => 9, 'symboltable' => {}, 'fatherblock' => {},
 'children' => [
 bless({
 'line' => 10, 'symboltable' => {}, 'fatherblock' => {},
 'children' => [
 bless({
 'line' => 10, 'symboltable' => {}, 'fatherblock' => {},
 'children' => [],
 }, 'BLOCK')
],
 }, 'BLOCK')
],
}, 'BLOCK')
],
}, 'BLOCK')
],
}, 'FUNCTION'),
bless({
 'function_name' => ['g', 13],

}, 'FUNCTION')
],
'line' => 1
}, 'PROGRAM');
$VAR1->{'children'}[0]['fatherblock'] = $VAR1;
$VAR1->{'children'}[0]['children'][0]['fatherblock'] = $VAR1->{'children'}[0];
$VAR1->{'children'}[0]['children'][0]['children'][0]['fatherblock']
 = $VAR1->{'children'}[0]['children'][0];
$VAR1->{'children'}[0]['children'][1]['fatherblock']
 = $VAR1->{'children'}[0];
$VAR1->{'children'}[0]['children'][1]['children'][0]['fatherblock']
 = $VAR1->{'children'}[0]['children'][1];
$VAR1->{'children'}[0]['children'][2]['fatherblock']
 = $VAR1->{'children'}[0];
$VAR1->{'children'}[0]['children'][2]['children'][0]['fatherblock']
 = $VAR1->{'children'}[0]['children'][2];
$VAR1->{'children'}[0]['children'][2]['children'][0]['children'][0]['fatherblock']
 = $VAR1->{'children'}[0]['children'][2]['children'][0];
$VAR1->{'children'}[1]['fatherblock'] = $VAR1;
$VAR1->{'children'}[1]['children'][0]['fatherblock'] = $VAR1->{'children'}[1];
$VAR1->{'children'}[1]['children'][1]['fatherblock'] = $VAR1->{'children'}[1];
$VAR1->{'children'}[1]['children'][1]['children'][0]['fatherblock']
 = $VAR1->{'children'}[1]['children'][1];
$VAR1->{'children'}[1]['children'][2]['fatherblock']
 = $VAR1->{'children'}[1];

```

Observe que este segundo código elimina las definiciones recursivas, retrasa las asignaciones y es código correcto. Es por tanto apto para reproducir la estructura de datos en un programa de prueba con, por ejemplo, `is_deeply`.

## Por que no Usar str en la Elaboración de las Pruebas



Puede que parezca una buena idea volcar los dos árboles esperado y obtenido mediante `str` y proceder a comparar las cadenas. Esta estrategia es inadecuado ya que depende de la versión de `Parse::Eyapp` con la que se esta trabajando. Si un usuario de nuestro módulo ejecuta la prueba con una versión distinta de `Parse::Eyapp` de la que hemos usado para la construcción de la prueba, puede obtener un "falso fallo" debido a que su version de `str` trabaja de forma ligeramente distinta.

### El Método `equal` de los Nodos

El método `equal` (version 1.094 de `Parse::Eyapp` o posterior) permite hacer comparaciones "difusas" de los nodos. El formato de llamada es:

```
$tree1->equal($tree2, attr1 => \&handler1, attr2 => \&handler2, ...)
```

Dos nodos se consideran iguales si:

1. Sus raíces `$tree1` y `$tree2` pertenecen a la misma clase
2. Tienen el mismo número de hijos
3. Para cada una de las claves especificadas `attr1`, `attr2`, etc. la existencia y definición es la misma en ambas raíces
4. Supuesto que en ambos nodos el atributo `attr` existe y está definido el manejador `handler($tree1, $tree2)` retorna cierto cuando es llamado
5. Los hijos respectivos de ambos nodos son iguales (en sentido recursivo o inductivo)

Sigue un ejemplo:

```
pl@nereida:~/LEyapp/examples$ cat -n equal.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Parse::Eyapp::Node;
4
5 my $string1 = shift || 'ASSIGN(VAR(TERMINAL))';
6 my $string2 = shift || 'ASSIGN(VAR(TERMINAL))';
7 my $t1 = Parse::Eyapp::Node->new($string1, sub { my $i = 0; $_->{n} = $i++ for @_ });
8 my $t2 = Parse::Eyapp::Node->new($string2);
9
10 # Without attributes
11 if ($t1->equal($t2)) {
12 print "\nNot considering attributes: Equal\n";
13 }
14 else {
15 print "\nNot considering attributes: Not Equal\n";
16 }
17
18 # Equality with attributes
19 if ($t1->equal($t2, n => sub { return $_[0] == $_[1] })) {
20 print "\nConsidering attributes: Equal\n";
21 }
22 else {
23 print "\nConsidering attributes: Not Equal\n";
24 }
```

## Usando equal en las Pruebas

Cuando desarrolle pruebas y desee obtener una comparación parcial del AST esperado con el AST obtenido puede usar la siguiente metodología:

1. Vuelque el árbol desde su programa de desarrollo con `Data::Dumper`
2. Compruebe que el árbol es correcto
3. Decida que atributos quiere comparar
4. Escriba la comparación de las estructuras y de los atributos pegando la salida de `Data::Dumper` y usando `equal` (versión 1.096 de `Parse::Eyapp`) como en el siguiente ejemplo:

```
pl@nereida:~/LEyapp/examples$ cat -n testequal.pl
```

```
1 #!/usr/bin/perl -w
2 use strict;
3 use Parse::Eyapp::Node;
4 use Data::Dumper;
5 use Data::Compare;
6
7 my $debugging = 0;
8
9 my $handler = sub {
10 print Dumper($_[0], $_[1]) if $debugging;
11 Compare($_[0], $_[1])
12 };
```

El manejador `$handler` es usado para comparar atributos (véanse las líneas 104-109). Copie el resultado de `Data::Dumper` en la variable `$t1`:

```
14 my $t1 = bless({
15 'types' => {
16 'CHAR' => bless({ 'children' => [] }, 'CHAR'),
17 'VOID' => bless({ 'children' => [] }, 'VOID'),
18 'INT' => bless({ 'children' => [] }, 'INT'),
19 'F(X_0(),INT)' => bless({
20 'children' => [
21 bless({ 'children' => [] }, 'X_0'),
22 bless({ 'children' => [] }, 'INT')]
23 }, 'F')
24 },
25 'symboltable' => { 'f' => { 'type' => 'F(X_0(),INT)', 'line' => 1 } },
26 'lines' => 2,
27 'children' => [
28 bless({
29 'symboltable' => {},
30 'fatherblock' => {},
31 'children' => [],
32 'depth' => 1,
33 'parameters' => [],
34 'function_name' => ['f', 1],
35 'symboltableLabel' => {},
36 'line' => 1
37 }, 'FUNCTION')
38],
```

```

39 'depth' => 0,
40 'line' => 1
41 }, 'PROGRAM');
42 $t1->{'children'}[0]['fatherblock'] = $t1;

```

Para ilustrar la técnica creamos dos árboles \$t2 y \$t3 similares al anterior que compararemos con \$t1. De hecho se han obtenido del texto de \$t1 suprimiendo algunos atributos. El árbol de \$t2 comparte con \$t1 los atributos "comprobados" mientras que no es así con \$t3:

```

44 # Tree similar to $t1 but without some attributes (line, depth, etc.)
45 my $t2 = bless({
46 'types' => {
47 'CHAR' => bless({ 'children' => [] }, 'CHAR'),
48 'VOID' => bless({ 'children' => [] }, 'VOID'),
49 'INT' => bless({ 'children' => [] }, 'INT'),
50 'F(X_0(),INT)' => bless({
51 'children' => [
52 bless({ 'children' => [] }, 'X_0'),
53 bless({ 'children' => [] }, 'INT')]
54 }, 'F')
55 },
56 'symboltable' => { 'f' => { 'type' => 'F(X_0(),INT)', 'line' => 1 } },
57 'children' => [
58 bless({
59 'symboltable' => {},
60 'fatherblock' => {},
61 'children' => [],
62 'parameters' => [],
63 'function_name' => ['f', 1],
64 }, 'FUNCTION')
65],
66 }, 'PROGRAM');
67 $t2->{'children'}[0]['fatherblock'] = $t2;
68
69 # Tree similar to $t1 but without some attributes (line, depth, etc.)
70 # and without the symboltable attribute
71 my $t3 = bless({
72 'types' => {
73 'CHAR' => bless({ 'children' => [] }, 'CHAR'),
74 'VOID' => bless({ 'children' => [] }, 'VOID'),
75 'INT' => bless({ 'children' => [] }, 'INT'),
76 'F(X_0(),INT)' => bless({
77 'children' => [
78 bless({ 'children' => [] }, 'X_0'),
79 bless({ 'children' => [] }, 'INT')]
80 }, 'F')
81 },
82 'children' => [
83 bless({
84 'symboltable' => {},
85 'fatherblock' => {},
86 'children' => [],
87 'parameters' => [],
88 'function_name' => ['f', 1],

```

```

89 }, 'FUNCTION')
90],
91 }, 'PROGRAM');
92
93 $t3->{'children'}[0]['fatherblock'] = $t2;

```

Ahora realizamos las comprobaciones de igualdad mediante `equal`:

```

95 # Without attributes
96 if (Parse::Eyapp::Node::equal($t1, $t2)) {
97 print "\nNot considering attributes: Equal\n";
98 }
99 else {
100 print "\nNot considering attributes: Not Equal\n";
101 }
102
103 # Equality with attributes
104 if (Parse::Eyapp::Node::equal(
105 $t1, $t2,
106 symboltable => $handler,
107 types => $handler,
108)
109) {
110 print "\nConsidering attributes: Equal\n";
111 }
112 else {
113 print "\nConsidering attributes: Not Equal\n";
114 }
115
116 # Equality with attributes
117 if (Parse::Eyapp::Node::equal(
118 $t1, $t3,
119 symboltable => $handler,
120 types => $handler,
121)
122) {
123 print "\nConsidering attributes: Equal\n";
124 }
125 else {
126 print "\nConsidering attributes: Not Equal\n";
127 }

```

Dado que los atributos usados son `symboltable` y `types`, los árboles `$t1` y `$t2` son considerados equivalentes. No así `$t1` y `$t3`. Observe el modo de llamada `Parse::Eyapp::Node::equal` como subrutina, no como método. Se hace así porque `$t1`, `$t2` y `$t3` no son objetos `Parse::Eyapp::Node`. La salida de `Data::Dumper` reconstruye la forma estructural de un objeto pero no reconstruye la información sobre la jerarquía de clases.

```
pl@nereida:~/LEyapp/examples$ testequal.pl
```

```
Not considering attributes: Equal
```

```
Considering attributes: Equal
```

```
Considering attributes: Not Equal
```

## 12.6. Análisis de Ámbito con Parse::Eyapp::Scope

Para calcular la relación entre una instancia de un identificador y su declaración usaremos el módulo `Parse::Eyapp::Scope` :

```
l@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '1,/~ *$/p' Scope.eypp | cat -n
1 /*
2 File: lib/Simple/Scope.eypp
3 Full Scope Analysis
4 Test it with:
5 lib/Simple/
6 eyapp -m Simple::Scope Scope.eypp
7 treereg -nonumbers -m Simple::Scope Trans
8 script/
9 usescope.pl prueba12.c
10 */
11 %{
12 use strict;
13 use Data::Dumper;
14 use List::MoreUtils qw(firstval lastval);
15 use Simple::Trans;
16 use Parse::Eyapp::Scope qw(:all);
```

Véamos el manejo de las reglas de `program`:

```
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/~^program:/,/~/~;$/p' Scope.eypp | cat
1 program:
2 {
3 reset_file_scope_vars();
4 }
5 definition<%name PROGRAM +>.program
6 {
7 $program->{symboltable} = { %st }; # creates a copy of the s.t.
8 for (keys %type) {
9 $type{$_} = Parse::Eyapp::Node->hnew($_);
10 }
11 $program->{depth} = 0;
12 $program->{line} = 1;
13 $program->{types} = { %type };
14 $program->{lines} = $tokenend;
15
16 my ($nondec, $declared) = $ids->end_scope($program->{symboltable}, $program, 'type'
17
18 if (@$nondec) {
19 warn "Identifier ".$_->key." not declared at line ".$_->line."\n" for @$nondec;
20 die "\n";
21 }
22
23 # Type checking: add a direct pointer to the data-structure
24 # describing the type
25 $_->{t} = $type{$_->{type}} for @$declared;
26
27 my $out_of_loops = $loops->end_scope($program);
28 if (@$out_of_loops) {
29 warn "Error: ".$ref($_). " outside of loop at line $_->{line}\n" for @$out_of_loops
```

```

30 die "\n";
31 }
32
33 # Check that are not dangling breaks
34 reset_file_scope_vars();
35
36 $program;
37 }
38 ;

```

Antes de comenzar la construcción del AST se inicializan las variables visibles en todo el fichero:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^sub reset_file_scope_vars/,/^}$/p'
1 sub reset_file_scope_vars {
2 %st = (); # reset symbol table
3 ($tokenbegin, $tokenend) = (1, 1);
4 %type = (INT => 1,
5 CHAR => 1,
6 VOID => 1,
7);
8 $depth = 0;
9 $ids = Parse::Eyapp::Scope->new(
10 SCOPE_NAME => 'block',
11 ENTRY_NAME => 'info',
12 SCOPE_DEPTH => 'depth',
13);
14 $loops = Parse::Eyapp::Scope->new(
15 SCOPE_NAME => 'exits',
16);
17 $ids->begin_scope();
18 $loops->begin_scope(); # just for checking
19 }

```

### El Método `Parse::Eyapp::Scope->new`

El método `Parse::Eyapp::Scope->new` crea un objeto del tipo manejador de ámbito.

Un *manejador de ámbito* es un objeto que ayuda en el cómputo de la función que asigna a cada *nodo instancia* de un nombre (variable, función, etiqueta, constante, identificador de tipo, etc.) su *declaración*, esto es su entrada en su tabla de símbolos.

En la línea 9 creamos el manejador de ámbito de los objetos identificadores `$ids` (ámbito de variables, funciones, etc.). En la línea 14 creamos un manejador de ámbito para los bucles (sentencias `CONTINUE`, `break`, etc.). El manejo de ámbito de los bucles consiste en asignar cada ocurrencia de una sentencia `BREAK` o `CONTINUE` con el bucle en el que ocurre. Por ejemplo:

```

pl@nereida:~/Lbook/code/Simple-Scope/script$ usescope.pl outbreak.c 2
1 test (int n, int m)
2 {
3 break;
4 while (n > 0) {
5 if (n>m) {
6 break;
7 }
8 else if (m>n){
9 continue;
10 }

```

```

11 n = n-1;
12 }
13 }

```

Error: BREAK outside of loop at line 3

En esta sección mostraremos como usar `Parse::Eyapp::Scope` cuando trabajemos en el análisis de condiciones dependientes del contexto.

La filosofía de `Parse::Eyapp::Scope` es que existen tres tipos de nodos en el AST:

1. Nodos que definen un ámbito y que tienen asociado un atributo 'tabla de símbolos', por ejemplo: Nodos programa, nodos función, nodos bloque conteniendo declaraciones, etc. Además del atributo tabla de símbolos es común que dichos nodos dispongan de un atributo 'profundidad de ámbito' que indique su nivel de anidamiento cuando el lenguaje siendo analizado usa ámbitos anidados.
2. Nodos que conllevan un uso de un nombre (*nodos de uso*): por ejemplo, nodos de uso de una variable en una expresión. El propósito del análisis de ámbito es dotar a cada uno de estos nodos de uso con un atributo 'scope' que referencia al *nodo ámbito* en el que se ha guardado la información que define las propiedades del objeto. Es posible que además queramos tener en dicho nodo un atributo 'entry' que sea una referencia directa a la entrada en la tabla de símbolos asociada con el nombre.
3. Otros tipos de nodo. Estos últimos pueden ser ignorados desde el punto de vista del análisis de ámbito

La asignación de ámbito se implanta a través de atributos que se añaden a las nodos de uso y a los nodos ámbito.

Algunos de los nombres de dichos atributos pueden ser especificados mediante los parámetros de `new`. En concreto:

- `SCOPE_NAME` es el nombre del atributo de la instancia que contendrá la referencia al nodo 'bloque' en el que ocurre esta instancia. Si no se especifica toma el valor `scope`.

En el ejemplo:

```

9 $ids = Parse::Eyapp::Scope->new(
10 SCOPE_NAME => 'block',
11 ENTRY_NAME => 'info',
12 SCOPE_DEPTH => 'depth',
13);

```

cada nodo asociado con una instancia de uso tendrá un atributo con clave `block` que será una referencia al nodo `BLOCK` en el que fué declarado el identificador.

- `ENTRY_NAME` es el nombre del atributo de la instancia que contendrá la referencia a la entrada de símbolos que corresponde a esta instancia. Si no se especifica tomará el valor `entry`.

En el ejemplo que estamos trabajando cada nodo `BLOCK` tendrá un atributo `symboltable` que será una referencia a la tabla de símbolos asociada con el bloque. Dado que el atributo `block` de un nodo `$n` asociado con una instancia de un identificador `$id` apunta al nodo `BLOCK` en el que se define, siempre sería posible acceder a la entrada de la tabla de símbolos mediante el código:

```
$n->{block}{symboltable}{$id}
```

El atributo `entry` crea una referencia directa en el nodo:

```
$n->{entry} = $n->{block}{symboltable}{$id}
```

de esta forma es mas directo acceder a la entrada de símbolos de una instancia de uso de un identificador.

Obviamente el atributo **entry** no tiene sentido en aquellos casos en que el análisis de ámbito no requiere de la presencia de tablas de símbolos, como es el caso del análisis de ámbito de las sentencias de cambio de flujo en bucles:

```
14 $loops = Parse::Eyapp::Scope->new(
15 SCOPE_NAME => 'exits',
16);
```

- **SCOPE\_DEPTH** es el nombre del atributo del nodo ámbito (nodo bloque) y contiene la profundidad de anidamiento del ámbito. Es opcional. Si no se especifica no será guardado.

### El Método **begin\_scope**

Este método debe ser llamado cada vez que se entra en una nueva región de ámbito. `Parse::Eyapp::Scope` asume un esquema de ámbitos léxicos anidados como ocurre en la mayoría de los lenguajes de programación: se supone que todos los nodos declarados mediante llamadas al método `scope_instance` entre dos llamadas consecutivas a `begin_scope` y `end_scope` definen la región del ámbito. Por ejemplo, en el análisis de ámbito de SimpleC llamamos a `begin_scope` cada vez que se entra en una definición de función:

```
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^funcDef:/,/~/;$/p' Scope.eypp | cat
1 funcDef:
2 $ID
3 {
4 $ids->begin_scope();
5 }
6 '(' $params ') '
7 $block
8 {
..
35 }
36 ;
```

y cada vez que se entra en un bloque:

```
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^block:/,/~/;$/p' Scope.eypp | cat -n
1 block:
2 '{'.bracket
3 { $ids->begin_scope(); }
4 declaration<%name DECLARATIONS *>.decs statement<%name STATEMENTS *>.sts '}'
5 {
25
26 }
27
28 ;
```

En el caso del análisis de ámbito de bucles la entrada en un bucle crea una nueva región de ámbito:

```
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^loopPrefix:/,/~/;$/p' Scope.eypp | cat -n
1 loopPrefix:
2 $WHILE '(' expression ')'
3 {
4 $loops->begin_scope;
5 $_[3]->{line} = $WHILE->[1];
```



```

6 $_[3]
7 }
8 ;

```

## El Método `end_scope`

En el código que sigue el método `end_scope` es llamado con tres argumentos:

```
my ($nodec, $dec) = $ids->end_scope($st, $block, 'type');
```

El significado de estos argumentos es el siguiente:

1. Una referencia a un hash `$st`. Este hash es la tabla de símbolos asociada con el bloque actual. Se asume que la clave de entrada de un nodo `$n` del árbol se obtiene mediante una llamada al método del nodo `key` definido por el programador: `$st->{$n->key}`. Se asume también que los valores del hash son una referencia a un hash conteniendo los atributos asociados con la clave `$n->key`.
2. Una referencia al nodo bloque `$block` o nodo de ámbito asociado con el hash. En nuestro ejemplo del análisis de ámbito de los identificadores las instancias de identificadores declarados en el presente bloque serán decorados con un atributo `block` que apunta a dicho nodo. El nombre es `block` porque así se especificó en la llamada a `new`:

```

9 $ids = Parse::Eyapp::Scope->new(
10 SCOPE_NAME => 'block',
11 ENTRY_NAME => 'info',
12 SCOPE_DEPTH => 'depth',
13);

```

3. Los argumentos adicionales como `'type'` que sean pasados a `end_scope` son interpretados como claves del hash apuntado por su entrada en la tabla de símbolos `$n->key`. Para cada uno de esos argumentos se crean referencias directas en el nodo `$n` a los mismos:

```
$n->{type} = $st->{$n->key}{type}
```

La llamada `$ids->end_scope($st, $block, 'type')` ocurre después del análisis de todo el bloque de la función:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^funcDef:/,/~/;$/p' Scope.eyp | cat
1 funcDef:
2 $ID
3 {
4 $ids->begin_scope();
5 }
6 '(' $params ')'
7 $block
8 {
9 my $st = $block->{symboltable};
10 my @decs = $params->children();
11 $block->{parameters} = [];
12 while (my ($bt, $id, $arrspec) = splice(@decs, 0, 3)) {
13
24 }
25 $block->{function_name} = $ID;
26 $block->type("FUNCTION");
27

```

```

28 my ($nodec, $dec) = $ids->end_scope($st, $block, 'type');
29
30 # Type checking: add a direct pointer to the data-structure
31 # describing the type
32 $_->{t} = $type{$_->{type}} for @$dec;
33
34 return $block;
35 }
36 ;

```

Todos los nodos `$n` del árbol que fueron declarados como instancias mediante llamadas al método `scope_instance` desde la última llamada a `begin_scope` son buscados en el hash referenciado por `$st`. Si la clave `$n->key` asociada con el nodo `$n` se encuentra entre las claves del hash (esto es, `exists $st->{$n->key}`) el nodo se anota como *declarado*. Los nodos-instancia para los cuales no existe una entrada en la tabla de símbolos se consideran *no declarados*.

Cuando un nodo se determina como declarado se establecen los atributos de declaración:

- Una referencia a su entrada en la tabla de símbolos (cuyo nombre fue establecido mediante el parámetro `ENTRY_NAME` de `new`)

```
$n->{$ids->{ENTRY_NAME}} = $st->{$n->key}
```

Si por ejemplo el nodo siendo usado se corresponde con el uso de una variable de nombre `a`, entonces `$n->key` deberá retornar `a` y la asignación anterior, asumiendo la declaración previa del manejador de ámbito `$ids`, sería:

```
$n->{info} = $st->{a}
```

- Una referencia al bloque en el que fué declarado (el nombre del atributo será el establecido en `SCOPE_NAME`).

```
$n->{$ids->{SCOPE_NAME}} = $block
```

Si por ejemplo el nodo siendo usado se corresponde con el uso de una variable de nombre `a`, la asignación anterior, asumiendo la declaración previa del manejador de ámbito `$ids`, sería:

```
$n->{block} = $block
```

donde `$block` es el nodo asociado con el bloque en el que fue declarado `$a`.

- Cualesquiera argumentos adicionales - como `'type'` en el ejemplo - que sean pasados a `end_scope` son interpretados como claves del hash apuntado por su entrada en la tabla de símbolos `$st->{$n->key}`. Para cada uno de esos argumentos se crean referencias directas en el nodo `$n` a los mismos.

```
$n->{type} = $st->{$n->key}{type}
```

esto es:

```
$n->{type} = $st->{a}{type}
```

Se supone que antes de la llamada a `end_scope` las declaraciones de los diferentes objetos han sido procesadas y la tabla hash ha sido rellenada con las mismas.

Los nodos `$n` que no aparezcan entre los declarados en este ámbito son apilados en la esperanza de que hayan sido declarados en un ámbito superior.

Para un ejemplo de uso véase la línea 28 en el código anterior. El tercer argumento **type** indica que para cada instancia de variable global que ocurre en el ámbito de **program** queremos que se cree una referencia desde el nodo a su entrada **type** en la tabla de símbolos. De este modo se puede conseguir un acceso mas rápido al tipo de la instancia si bien a costa de aumentar el consumo de memoria.

En un contexto de lista **end\_scope** devuelve un par de referencias. La primera es una referencia a la lista de instancias/nodos del ámbito actual que no aparecen como claves del hash. La segunda a los que aparecen.

En un contexto escalar devuelve la lista de instancias/nodos no declarados.

## El Método **key** para Nodos con Ámbito

El programador deberá proveer a cada clase de nodo (subclases de **Parse::Eyapp::Node**) que pueda ser instanciado/usado en un ámbito o bloque de un método **key**. En nuestro caso los nodos del tipo **VAR**, **VARARRAY** y **FUNCTIONCALL** son los que pueden ser usados dentro de expresiones y sentencias.

El método **key** se usa para computar el valor de la clave de entrada en el hash para esa clase de nodo.

Para la tabla de símbolos de los identificadores en SimpleC necesitamos definir el método **key** para los nodos **VAR**, **VARARRAY** y **FUNCTIONCALL**:

```
827 sub VAR::key {
828 my $self = shift;
829
830 return $self->child(0)->{attr}[0];
831 }
832
833 *VARARRAY::key = *FUNCTIONCALL::key = \&VAR::key;
```

Si tiene dudas repase la definición de **Variable** en la descripción de la gramática en la página 640:

```
Variable:
 %name VAR
 ID
| %name VARARRAY
 %ID ('[' binary ']') <%name INDEXSPEC +>
;
```

y el subárbol para la asignación de **a = 2** en la página 648.

```
ASSIGN(
 VAR(
 TERMINAL[a:5]
),
 INUM(
 TERMINAL[2:5]
)
) # ASSIGN,
```

## El Modo de Llamada Simplificado a **end\_scope**

En el caso de un análisis de ámbito simple como es el análisis de ámbito de los bucles no se requiere de la presencia de una tabla de símbolos. El único argumento de **end\_scope** es la referencia al nodo que define el ámbito.

```
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^statement:\/,\/^;$\/p' Scope.eyp | ca
1 statement:
2 expression ';' { $_[1] }
```

```

3 | ';'
4 | %name BREAK
5 $BREAK ';'
6 {
7 my $self = shift;
8 my $node = $self->YYBuildAST(@_);
9 $node->{line} = $BREAK->[1];
10 $loops->scope_instance($node);
11 return $node;
12 }
13 | %name CONTINUE
14 $CONTINUE ';'
15 {
16
21 }
..
34 | %name WHILE
35 $loopPrefix statement
36 {
37 my $self = shift;
38 my $wnode = $self->YYBuildAST(@_);
39 $wnode->{line} = $loopPrefix->{line};
40 my $breaks = $loops->end_scope($wnode);
41 return $wnode;
42 }
43 ;

```

### El Método `scope_instance`

Este método del objeto `Parse::Eyapp::Scope` empuja el nodo que se le pasa como argumento en la cola de instancias del manejador de ámbito. El nodo se considerará una ocurrencia de un objeto dentro del ámbito actual (el que comienza en la última ejecución de `begin_scope`). En el compilador de SimpleC debemos hacer:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^Primary:/,/^$;/p' Scope.eyyp | cat
1 Primary:
2 %name INUM
3 INUM
4 | %name CHARCONSTANT
5 CHARCONSTANT
6 | $Variable
7 {
8 $ids->scope_instance($Variable);
9 return $Variable
10 }
11 | '(' expression ')' { $_[2] }
12 | $function_call
13 {
14 $ids->scope_instance($function_call);
15 return $function_call # bypass
16 }
17 ;

```

Otros lugares en los que ocurren instancias de identificadores son las asignaciones:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^binary:/,/^$;/p' Scope.eypp | cat -
1 binary:
2 Unary { $_[1] }
3 | %name PLUS
4 binary '+' binary
..
23 | %name ASSIGN
24 $Variable '=' binary
25 {
26 goto &declare_instance_and_build_node;
27 }
28 | %name PLUSASSIGN
29 $Variable '+=' binary
30 {
31 goto &declare_instance_and_build_node;
32 }
33 | %name MINUSASSIGN
34 $Variable '-=' binary
35 {
36 goto &declare_instance_and_build_node;
37 }
38 | %name TIMESASSIGN
39 $Variable '*=' binary
40 {
41 goto &declare_instance_and_build_node;
42 }
43 | %name DIVASSIGN
44 $Variable '/=' binary
45 {
46 goto &declare_instance_and_build_node;
47 }
48 | %name MODASSIGN
49 $Variable '%=' binary
50 {
51 goto &declare_instance_and_build_node;
52 }
53 ;

```

Como indica su nombre, la función `declare_instance_and_build_node` declara la instancia y crea el nodo del AST:

```

116 sub declare_instance_and_build_node {
117 my ($parser, $Variable) = @_[0,1];
118
119 $ids->scope_instance($Variable);
120 goto &Parse::Eyapp::Driver::YYBuildAST;
121 }

```

En el caso del análisis de ámbito en bucles las instancias ocurren en las sentencias de `break` y `continue`:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^statement:/,/^$;/p' Scope.eypp | cat -
1 statement:
2 expression ';' { $_[1] }
3 | ';'

```

```

4 | %name BREAK
5 $BREAK ';'
6 {
7 my $self = shift;
8 my $node = $self->YYBuildAST(@_);
9 $node->{line} = $BREAK->[1];
10 $loops->scope_instance($node);
11 return $node;
12 }
13 | %name CONTINUE
14 $CONTINUE ';'
15 {
16 my $self = shift;
17 my $node = $self->YYBuildAST(@_);
18 $node->{line} = $CONTINUE->[1];
19 $loops->scope_instance($node);
20 return $node;
21 }
..
34 | %name WHILE
35 $loopPrefix statement
36 {
37 my $self = shift;
38 my $wnode = $self->YYBuildAST(@_);
39 $wnode->{line} = $loopPrefix->{line};
40 my $breaks = $loops->end_scope($wnode);
41 return $wnode;
42 }
43 ;

```

**Cálculo del Ámbito en Bloques** El modo de uso se ilustra en el manejo de los bloques

```
block: '{' declaration * statement * '}'
```

La computación del ámbito requiere las siguientes etapas:

1. Es necesario introducir una acción intermedia (línea 3) para indicar que una "llave abrir" determina el comienzo de un bloque.
2. Durante las sucesivas visitas a **declaration** construimos tablas de símbolos que son mezcladas en las líneas 8-15.
3. Todas las instancias de nodos-nombre que ocurren cuando se visitan los hijos de **statements** son declaradas como instanciaciones con **scope\_instance** (véanse los acciones semánticas para **Primary** en la página 677).
4. La llamada a **end\_scope** de la línea 19 produce la computación parcial de la función de asignación de ámbito para los nodos/nombre que fueron instanciados en este bloque.

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/^block:\/,\/~;$$/p' Scope.eypp | cat -n
1 block:
2 '{' .bracket
3 { $ids->begin_scope(); }
4 declaration<%name DECLARATIONS *>.decs statement<%name STATEMENTS *>.sts '}'
5 {
6 my %st;

```

```

7
8 for my $lst ($decs->children) {
9
10 # control duplicated declarations
11 my $message;
12 die $message if $message = is_duplicated(\%st, $lst);
13
14 %st = (%st, %$lst);
15 }
16 $sts->{symboltable} = \%st;
17 $sts->{line} = $bracket->[1];
18 $sts->type("BLOCK") if (%st);
19 my ($nondec, $dec) = $ids->end_scope(\%st, $sts, 'type');
20
21 # Type checking: add a direct pointer to the data-structure
22 # describing the type
23 $_->{t} = $type{$_->{type}} for @$dec;
24
25 return $sts;
26 }
27
28 ;

```

Actualmente las acciones intermedias tal y como la que se ven en el código anterior se usan habitualmente para producir efectos laterales en la construcción del árbol. No dan lugar a la inserción de un nodo. Esto es, la variable auxiliar/temporal generada para dar lugar a la acción intermedia no es introducida en el árbol generado.

## La Jerarquía de Bloques

Para completar el análisis de ámbito queremos decorar cada nodo **BLOCK** con un atributo **fatherblock** que referencia al nodo bloque que lo contiene. Para esta fase usaremos el lenguaje **Treeregexp**. Las frases de este lenguaje permiten definir conjuntos de árboles. Cuando van en un fichero aparte los programas **Treeregexp** suelen tener la extensión **trg**. Por ejemplo, el siguiente programa

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ cat -n Trans.trg
1 /* Scope Analysis */
2 blocks: /BLOCK|FUNCTION|PROGRAM/
3
4 retscope: /FUNCTION|RETURN/
5
6 loop_control: /BREAK|CONTINUE|WHILE/

```

define tres expresiones árbol a las que nos podremos referir en el cliente mediante las variables **\$blocks**, **\$retscope** y **loops\_control**. La primera expresión árbol tiene por nombre **blocks** y casará con cualesquiera árboles de las clases **BLOCK**, **FUNCTION** o **PROGRAM**. Es necesario compilar el programa **Treeregexp**:

```

l@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ treereg -m Simple::Scope Trans
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ ls -ltr | tail -2
-rw-r--r-- 1 pl users 122 2007-12-10 11:49 Trans.trg
-rw-r--r-- 1 pl users 1544 2007-12-10 11:49 Trans.pm

```

El módulo generado **Trans.pm** contendrá una subrutina por cada expresión regular árbol en el programa **Treeregexp**. La subrutina devuelve cierto si el nodo que se le pasa como primer argumento casa con la expresión árbol. La subrutina espera como primer argumento el nodo, como segundo argumento el padre del nodo y como tercero el objeto **Parse::Yapp::YATW** que encapsula la transformación:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '19,31p' Trans.pm
sub blocks {
 my $blocks = $_[3]; # reference to the YATW pattern object
 my $W;

 {
 my $child_index = 0;

 return 0 unless ref($W = $_[$child_index]) =~ m{\bBLOCK\b|\bFUNCTION\b|\bPROGRAM\b}x;

 } # end block of child_index
 1;

} # end of blocks

```

No solo se construye una subrutina `blocks` sino que en el espacio de nombres correspondiente se crean objetos `Parse::Eyapp::YATW` por cada una de las transformaciones especificadas en el programa `trg`, como muestra el siguiente fragmento del módulo conteniendo el código generado por `treeregexp`:

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '16p' Trans.pm
our @all = (our $blocks, our $retscope, our $loop_control,) =
 Parse::Eyapp::YATW->buildpatterns(
 blocks => \&blocks,
 retscope => \&retscope,
 loop_control => \&loop_control,
);

```

Estos objetos `Parse::Eyapp::YATW` estan disponibles en el programa `eyapp` ya que hemos importado el módulo:

```

1@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '1,/^\s*$/p' Scope.eyp | cat -n
1 /*
2 File: lib/Simple/Scope.eyp
3 Full Scope Analysis
4 Test it with:
5 lib/Simple/
6 eyapp -m Simple::Scope Scope.eyp
7 treereg -nonumbers -m Simple::Scope Trans
8 script/
9 usescope.pl prueba12.c
10 */
11 %{
12 use strict;
13 use Data::Dumper;
14 use List::MoreUtils qw(firstval lastval);
15 use Simple::Trans;
16 use Parse::Eyapp::Scope qw(:all);

```

Es por eso que, una vez realizadas las fases de análisis léxico, sintáctico y de ámbito podemos usar el objeto `$blocks` y el método `m` (por `matching`) de dicho objeto (véase la línea 14 en el código de `compile`):

```

pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '/sub compile/,/^\s*$/p' Scope.eyp | cat -n
1 sub compile {
2 my($self)=shift;
3

```



```

4 my ($t);
5
6 $self->YYData->{INPUT} = shift;
7
8 $t = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
9 #yydebug => 0x1F
10);
11
12 # Scope Analysis: Block Hierarchy
13 our $blocks;
14 my @blocks = $blocks->m($t);
15 $_->node->{fatherblock} = $_->father->{node} for (@blocks[1..$#blocks]);
16
17
28
29 return $t;
30 }

```

El método `m` para árboles y expresiones regulares árbol YATW funciona de manera parecida al método `m` para cadenas y expresiones regulares convencionales. El resultado de un matching en árboles es un árbol. Si se trabaja en un contexto de lista es una lista de árboles.

La llamada `$blocks->m($t)` permite la búsqueda de los nodos de `$t` que casan con la expresión regular árbol para `blocks`. En un contexto de lista `m` devuelve una lista con nodos del tipo `Parse::Eyapp::Node::Match` que referencian a los nodos que han casado. Los nodos `Parse::Eyapp::Node::Match` son a su vez nodos (heredan de) `Parse::Eyapp::Node`. Los nodos aparecen en la lista retornada en orden primero profundo de recorrido del árbol `$t`.

Los nodos en la lista se estructuran según un árbol (atributos `children` y `father`) de manera que el padre de un nodo `$n` del tipo `Parse::Eyapp::Node::Match` es el nodo `$f` que referencia al inmediato antecesor en el árbol que ha casado.

Un nodo `$r` de la clase `Parse::Eyapp::Node::Match` dispone de varios atributos y métodos:

- El método `$r->node` retorna el nudo del árbol `$t` que ha casado
- El método `$r->father` retorna el nodo padre en el árbol the matching. Se cumple que `$r->father->node` es una referencia al antepasado mas cercano en `$t` de `$r->node` que casa con el patrón árbol.
- El método `$r->coord` retorna las coordenadas del nudo del árbol que ha casado usando una notación con puntos. Por ejemplo la coordenada `".1.3.2"` denota al nodo `$t->child(1)->child(3)->child(2)`, siendo `$t` la raíz del árbol de búsqueda.
- El método `$r->depth` retorna la profundidad del nudo del árbol que ha casado

En un contexto escalar `m` devuelve el primer elemento de la lista de nodos `Parse::Eyapp::Node::Match`.

```
pl@nereida:~/Lbook/code/Simple-Scope/script$ perl -wd usescope.pl blocks.c 2
```

```
Loading DB routines from perl5db.pl version 1.28
Editor support available.
```

```
Enter h or 'h h' for help, or 'man perldebug' for more help.
```

```
main:.(usescope.pl:6): my $filename = shift || die "Usage:\n$0 file.c\n";
DB<1> c Simple::Scope::compile
```

```
1 test (int n)
2 {
```

```

3 int a;
4
5 while (1) {
6 if (1>0) {
7 int a;
8 break;
9 }
10 else if (2>0){
11 char b;
12 continue;
13 }
14 }
15 }
Simple::Scope::compile(Scope.eyp:784):
784: my($self)=shift;
DB<2> 1 783,797
783 sub compile {
784==> my($self)=shift;
785
786: my ($t);
787
788: $self->YYData->{INPUT} = shift;
789
790: $t = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
791 #yydebug => 0x1F
792);
793
794 # Scope Analysis: Block Hierarchy
795: our $blocks;
796: my @blocks = $blocks->m($t);
797: $_->node->{fatherblock} = $_->father->{node} for (@blocks[1..$#blocks]);
DB<3> c 797
Simple::Scope::compile(Scope.eyp:797):
797: $_->node->{fatherblock} = $_->father->{node} for (@blocks[1..$#blocks]);
DB<4> x Parse::Eyapp::Node->str(@blocks)
0 '
Match[[PROGRAM:0:blocks]](
 Match[[FUNCTION:1:blocks:test]](
 Match[[BLOCK:6:blocks:10:4]],
 Match[[BLOCK:5:blocks:6:4]]
)
)'
1 '
Match[[FUNCTION:1:blocks:test]](
 Match[[BLOCK:6:blocks:10:4]],
 Match[[BLOCK:5:blocks:6:4]]
)'
2 '
Match[[BLOCK:5:blocks:6:4]]'
3 '
Match[[BLOCK:6:blocks:10:4]]'

```

La información que aparece en los nodos `Match` es como sigue:

- La clase del nodo al que referencia `PROGRAM`, `FUNCTION`, etc.

- La profundidad del nodo que ha casado en \$t
- Los nombres de las transformaciones con las que casó el nodo
- Si el nodo tiene un método `info` la información asociada. En este caso la línea y la profundidad en los bloques.

```
DB<5> p $t->str
```

```
PROGRAM^{0}(
 FUNCTION[test]^{1}(
 WHILE(
 INUM(TERMINAL[1:5]),
 STATEMENTS(
 IFELSE(
 GT(INUM(TERMINAL[1:6]), INUM(TERMINAL[0:6])),
 BLOCK[6:4]^{2}(
 BREAK
),
 IF(
 GT(INUM(TERMINAL[2:10]), INUM(TERMINAL[0:10])),
 BLOCK[10:4]^{3}(
 CONTINUE
)
)
)
)
)
)
)
....
DB<6> x map {$_->coord} @blocks
0 ''
1 '.0'
2 '.0.0.1.0.1'
3 '.0.0.1.0.2.1'
DB<7> p $t->descendant('.0.0.1.0.2.1')->str

BLOCK[10:4]^{0}(
 CONTINUE
)

0)
Symbol Table of block at line 10
$VAR1 = {
 'b' => {
 'type' => 'CHAR',
 'line' => 11
 }
};
DB<8> x map {$_->depth} @blocks
0 0
1 1
2 5
```

```

3 6
DB<9> x map {ref($_->node) } @blocks
0 'PROGRAM'
1 'FUNCTION'
2 'BLOCK'
3 'BLOCK'
DB<10> x map {ref($_->father) } @blocks
0 ''
1 'Parse::Eyapp::Node::Match'
2 'Parse::Eyapp::Node::Match'
3 'Parse::Eyapp::Node::Match'
DB<11> x map {ref($_->father->node) } @blocks[1..3]
0 'PROGRAM'
1 'FUNCTION'
2 'FUNCTION'
DB<12> x $blocks[2]->father->node->{function_name}
0 ARRAY(0x86ed84c)
0 'test'
1 1

```

## El Método names

La clase `Parse::Eyapp::Node::Match` dispone además de otros métodos para el caso en que se usan varios patrones en la misma búsqueda. Por ejemplo, el método `$r->names` retorna una referencia a los nombres de los patrones que han casado con el nodo. En el ejemplo que sigue aparecen tres transformaciones `fold`, `zxw` y `wxz`. La llamada a `m` de la línea 19 toma una forma diferente. ahora `m` es usado como método del árbol `$t` y recibe como argumentos la lista con los objetos expresiones regulares árbol (`Parse::Eyapp::YATW`).

```

pl@nereida:~/LEyapp/examples$ cat -n m2.pl
1 #!/usr/bin/perl -w
2 use strict;
3 use Rule6;
4 use Parse::Eyapp::Treeregexp;
5
6 Parse::Eyapp::Treeregexp->new(STRING => q{
7 fold: /TIMES|PLUS|DIV|MINUS/(NUM, NUM)
8 zxw: TIMES(NUM($x), .) and { $x->{attr} == 0 }
9 wxz: TIMES(., NUM($x)) and { $x->{attr} == 0 }
10 }->generate();
11
12 # Syntax analysis
13 my $parser = new Rule6();
14 my $input = "0*0*0";
15 my $t = $parser->Run(\$input);
16 print "Tree:", $t->str, "\n";
17
18 # Search
19 my $m = $t->m(our ($fold, $zxw, $wxz));
20 print "Match Node:\n", $m->str, "\n";

```

La primera expresión regular árbol `fold` casa con cualquier árbol tal que la clase del nodo raíz case con la expresión regular (clásica) `/TIMES|PLUS|DIV|MINUS/` y tenga dos hijos de clase `NUM`.

La segunda expresión regular `zxw` (por *zero times whatever*) casa con cualquier árbol tal que la clase del nodo raíz es `TIMES` y cuyo primer hijo pertenece a la clase `NUM`. Se debe cumplir además

que el hijo del nodo NUM (el nodo TERMINAL proveído por el analizador léxico) debe tener su atributo `attr` a cero. La notación NUM(\$x) hace que en \$x se almacene una referencia al nodo hijo de NUM. La expresión regular árbol wxz es la simétrica de zwx.

Cuando se ejecuta, el programa produce la salida:

```
pl@nereida:~/LEyapp/examples$ m2.pl
Tree:TIMES(TIMES(NUM(TERMINAL),NUM(TERMINAL)),NUM(TERMINAL))
Match Node:
Match[[TIMES:0:wxz]](Match[[TIMES:1:fold,zxw,wxz]])
```

## 12.7. Resultado del Análisis de Ámbito

Como se ha comentado en la sección 12.5 el volcado de los AST con `Data::Dumper` es insuficiente. Para la visualización utilizaremos el método `str` presentado en la sección 12.8.

### Código de Apoyo a la Visualización

Sigue el código de apoyo a la visualización:

```
pl@nereida:~/Lbook/code/Simple-Scope/lib/Simple$ sed -ne '845,$p' Scope.eypp | cat -n
1 ##### Debugging and Display
2 sub show_trees {
3 my ($t) = @_ ;
4
5 print Dumper $t if $debug > 1;
6 $Parse::Eyapp::Node::INDENT = 2;
7 $Data::Dumper::Indent = 1;
8 print $t->str."\n";
9 }
10
11 $Parse::Eyapp::Node::INDENT = 1;
12 sub TERMINAL::info {
13 my @a = join ':', @{$_[0]->{attr}};
14 return "@a"
15 }
16
17 sub PROGRAM::footnote {
18 return "Types:\n"
19 .Dumper($_[0]->{types}).
20 "Symbol Table of the Main Program:\n"
21 .Dumper($_[0]->{symboltable})
22 }
23
24 sub FUNCTION::info {
25 return $_[0]->{function_name}[0]
26 }
27
28 sub FUNCTION::footnote {
29 my $text = '';
30 $text .= "Symbol Table of $_[0]->{function_name}[0]\n" if $_[0]->type eq 'FUNCTION';
31 $text .= "Symbol Table of block at line $_[0]->{line}\n" if $_[0]->type eq 'BLOCK';
32 $text .= Dumper($_[0]->{symboltable});
33 return $text;
34 }
35
```

```

36 sub BLOCK::info {
37 my $info = "$_[0]->{line}:$_[0]->{depth}";
38 my $fatherblock = $_[0]->{fatherblock};
39 $info .= ":".$fatherblock->info
40 if defined($fatherblock) and UNIVERSAL::can($fatherblock, 'info');
41 return $info;
42 }
43
44 *BLOCK::footnote = \&FUNCTION::footnote;
45
46 sub VAR::info {
47 return $_[0]->{type} if defined $_[0]->{type};
48 return "No declarado!";
49 }
50
51 *FUNCTIONCALL::info = *VARARRAY::info = \&VAR::info;

```

### Un Ejemplo con Errores de Ámbito

Utilizando este código obtenemos una descripción completa y concisa del árbol anotado. Si hay errores de ámbito se mostrarán.

```
pl@nereida:~/Lbook/code/Simple-Scope/script$ usescope.pl prueba27.c 2
```

```

1 int a,b;
2
3 int f(char c) {
4 a[2] = 4;
5 {
6 int d;
7 d = a + b;
8 }
9 c = d * 2;
10 return g(c);
11 }

```

Identifier d not declared at line 9

Identifier g not declared at line 10

### Mostrando el Resultado del Análisis de Ámbito

Si no hay errores se obtiene un listado enumerado del programa fuente, el árbol, la tabla de tipos y las tablas de símbolos:

```
pl@nereida:~/Lbook/code/Simple-Scope/script$ usescope.pl prueba25.c 2
```

```

1 int a[20],b,e[10];
2
3 g() {}
4
5 int f(char c) {
6 char d;
7 c = 'X';
8 e[b] = 'A'+c;
9 {
10 int d;
11 d = d + b;

```

```

12 }
13 c = d * 2;
14 return c;
15 }
16
17
PROGRAM^{0}(
 FUNCTION[g]^{1},
 FUNCTION[f]^{2}(
 ASSIGN(
 VAR[CHAR] (
 TERMINAL[c:7]
),
 CHARCONSTANT(
 TERMINAL['X':7]
)
) # ASSIGN,
 ASSIGN(
 VARARRAY[A_10(INT)] (
 TERMINAL[e:8],
 INDEXSPEC(
 VAR[INT] (
 TERMINAL[b:8]
)
) # INDEXSPEC
) # VARARRAY,
 PLUS(
 CHARCONSTANT(
 TERMINAL['A':8]
),
 VAR[CHAR] (
 TERMINAL[c:8]
)
) # PLUS
) # ASSIGN,
 BLOCK[9:3:f]^{3}(
 ASSIGN(
 VAR[INT] (
 TERMINAL[d:11]
),
 PLUS(
 VAR[INT] (
 TERMINAL[d:11]
),
 VAR[INT] (
 TERMINAL[b:11]
)
) # PLUS
) # ASSIGN
) # BLOCK,
 ASSIGN(
 VAR[CHAR] (
 TERMINAL[c:13]

```

```

),
 TIMES(
 VAR[CHAR](
 TERMINAL[d:13]
),
 INUM(
 TERMINAL[2:13]
)
) # TIMES
) # ASSIGN,
RETURN(
 VAR[CHAR](
 TERMINAL[c:14]
)
) # RETURN
) # FUNCTION
) # PROGRAM

0)
Types:
$VAR1 = {
 'A_10(INT)' => bless({
 'children' => [
 bless({
 'children' => []
 }, 'INT')
]
 }, 'A_10'),
 'F(X_1(CHAR),INT)' => bless({
 'children' => [
 bless({
 'children' => [
 bless({
 'children' => []
 }, 'CHAR')
]
 }, 'X_1'),
 $VAR1->{'A_10(INT)'}{'children'}[0]
]
 }, 'F'),
 'CHAR' => $VAR1->{'F(X_1(CHAR),INT)'}{'children'}[0]{'children'}[0],
 'VOID' => bless({
 'children' => []
 }, 'VOID'),
 'INT' => $VAR1->{'A_10(INT)'}{'children'}[0],
 'A_20(INT)' => bless({
 'children' => [
 $VAR1->{'A_10(INT)'}{'children'}[0]
]
 }, 'A_20'),
 'F(X_0(),INT)' => bless({
 'children' => [
 bless({

```



```

 'children' => []
 }, 'X_0'),
 $VAR1->{'A_10(INT)'}{'children'}[0]
]
}, 'F')
};

```

Symbol Table of the Main Program:

```

$VAR1 = {
 'e' => {
 'type' => 'A_10(INT)',
 'line' => 1
 },
 'a' => {
 'type' => 'A_20(INT)',
 'line' => 1
 },
 'b' => {
 'type' => 'INT',
 'line' => 1
 },
 'g' => {
 'type' => 'F(X_0(),INT)',
 'line' => 3
 },
 'f' => {
 'type' => 'F(X_1(Char),INT)',
 'line' => 5
 }
};

```

---

1)  
Symbol Table of g  
\$VAR1 = {};

---

2)  
Symbol Table of f  
\$VAR1 = {  
 'c' => {  
 'type' => 'CHAR',  
 'param' => 1,  
 'line' => 5  
 },  
 'd' => {  
 'type' => 'CHAR',  
 'line' => 6  
 }  
};

---

3)  
Symbol Table of block at line 9

```
$VAR1 = {
 'd' => {
 'type' => 'INT',
 'line' => 10
 }
};
```

```
pl@nereida:~/Lbook/code/Simple-Scope/script$
```

Observe como se anotan las referencias a pie de árbol usando la notación `~{#}`. La salida se completa con las *notas a pie de árbol*.

## 12.8. Usando el Método `str` para Analizar el Árbol

### El Método `str` de los Nodos

Para dar soporte al análisis de los árboles y su representación, `Parse::Eyapp` provee el método `str`. El siguiente ejemplo con el depurador muestra el uso del método `str`. El método `Parse::Eyapp::Node::str` retorna una cadena que describe como término el árbol enraizado en el nodo que se ha pasado como argumento.

### El Método `info`

El método `str` cuando visita un nodo comprueba la existencia de un método `info` para la clase del nodo. Si es así el método será llamado. Obsérvese como en las líneas 5 y 6 proveemos de métodos `info` a los nodos `TERMINAL` y `FUNCTION`. La consecuencia es que en la llamada de la línea 7 el término que describe al árbol es decorado con los nombres de funciones y los atributos y números de línea de los terminales.

Los nodos de la clase `Parse::Eyapp::Node::Match` vienen con un método `info` asociado el cual muestra el tipo del nodo apuntado, su profundidad y el nombre del patrón que casó (línea 9 de la sesión).

Veamos un ejemplo que ilustra también algunos "trucos" de depuración. Se han introducido comentarios en el texto de salida:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> perl -wd usesimple4.pl
```

la opción `-w` por `warning` y `-d` para activar el debugger

```
Loading DB routines from perl5db.pl version 1.28
Editor support available.
```

El mensaje `Editor support available` nos señala la presencia de los servicios de edición de la entrada. Si no sale el mensaje instale `Term::ReadLine::Gnu` en su máquina!

```
main:.(usesimple4.pl:6): $Data::Dumper::Indent = 1;
```

La información entre paréntesis indica el fichero en que estamos y la línea (`usesimple4.pl`)

```
DB<1> f Simple4.eyp # <-- Quiero poner un break en Simple4.eyp,
```

El comando `f` es esencial cuando se usa `eyapp`. Le indicamos al debugger que queremos usar como fichero por defecto el de nuestra gramática. Los siguientes comandos se referirán a `Simple4.eyp`:

```
DB<2> b 654 # <--- break en la línea 654 de Simple4.eyp
DB<3> l 654,656 # <-- Se trata de las líneas en las que se establece
 # la jerarquía de bloques
654:b my @blocks = $SimpleTrans::blocks->m($t);
655: $_->node->{fatherblock} = $_->father->{node} \
 for (@blocks[1..$#blocks]);
```

```

656: $Data::Dumper::Deepcopy = 1;
DB<4> c # c indica que se continúa la ejecución del programa
 # hasta encontrar un "break"

```

La ejecución del programa continúa produciendo la salida del fuente hasta encontrar el punto de break:

```

char d,e[1][2];
char f() {
 int c[2];
 char d;

 return d;
}
Simple4::Run(Simple4.eyp:654): my @blocks = $SimpleTrans::blocks->m($t);
DB<4> n # <--- next ejecuta la siguiente instrucción. Si es una sub no se entra.
Simple4::Run(Simple4.eyp:655): $_->node->{fatherblock} = $_->father->{node} \
 for (@blocks[1..$#blocks]);
DB<4> # La simple pulsación de un retorno de carro
 # hace que se repita la última instrucción s o n
Simple4::Run(Simple4.eyp:655): $_->node->{fatherblock} = $_->father->{node} \
 for (@blocks[1..$#blocks]);
DB<4>
Simple4::Run(Simple4.eyp:656): $Data::Dumper::Deepcopy = 1;
DB<4>
Simple4::Run(Simple4.eyp:657): print Dumper $t;
DB<4> x $t->str # El método str para los objetos Parse::Eyapp::Node devuelve
 # una descripción del árbol
0 'PROGRAM(FUNCTION(RETURN(TERMINAL,VAR(TERMINAL))))'
DB<5> sub TERMINAL::info { @a = join ':', @{$_[0]->{attr}}; "[@a]"
DB<6> sub FUNCTION::info { "[".$_[0]->{function_name}[0]."]" }
DB<7> x $t->str
0 'PROGRAM(FUNCTION[f](RETURN(VAR(TERMINAL[d:6]))))'
 # La variable @Parse::Eyapp::Node::PREFIXES Controla que
 # prefijos se deben eliminar en la presentación
DB<8> p @Parse::Eyapp::Node::PREFIXES
Parse::Eyapp::Node::
DB<9> x $blocks[0]->str
0 'Match[PROGRAM:0:blocks](Match[FUNCTION:1:blocks:[f]])'
 # La información de un nodo Match es la
 # terna [type:depth:nombre de la treeregexp]

```

## Modo de uso de str

El método `str` ha sido concebido como una herramienta de ayuda a la depuración y verificación de los programas árbol. La metodología consiste en incorporar las pequeñas funciones `info` al programa en el que trabajamos. Así en `Simple4.eyp` añadimos el siguiente código:

```

640 sub Run {
641 my($self)=shift;
642
663 }
664
665 sub TERMINAL::info {
666 my @a = join ':', @{$_[0]->{attr}};

```

```

667 return "[@a]"
668 }
669
670 sub FUNCTION::info {
671 return "[".$_[0]->{function_name}[0]."]"
672 }
673
674 sub BLOCK::info {
675 return "[".$_[0]->{line}."]"
676 }

```

### Un ejemplo de Salida con Bloques Anidados

En el código mostrado no se ha realizado la optimización consistente en no incluir en el árbol de bloques a aquellos bloques cuya sección de declaraciones es vacía. Así podemos estudiar la estructura de datos resultante para un pequeño programa como el que sigue a continuación.

Cuadro 12.1: Representación de AST con str

| Programa                                                                                                                                              | \$t->str                                                                                                                                                                                                                                                                          | \$blocks[0]->str                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|-------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> 1 char d0; 2 char f() { 3   { 4     {} 5   } 6   { 7     { } 8   } 9   { 10    {} 11  } 12 } 13 g() { 14 {} 15 { 16   {} 17 } 18 {} 19 } </pre> | <pre> PROGRAM(   FUNCTION[f](     BLOCK[3](       BLOCK[4]     ),     BLOCK[6](       BLOCK[7]),     BLOCK[9](       BLOCK[10](         BLOCK[10]       )     ),     FUNCTION[g](       BLOCK[14],       BLOCK[15](         BLOCK[16]       ),       BLOCK[18]     )   ) ) </pre> | <pre> Match[PROGRAM:0:blocks](   Match[FUNCTION:1:blocks:[f]](     Match[BLOCK:2:blocks:[3]](       Match[BLOCK:3:blocks:[4]]     ),     Match[BLOCK:2:blocks:[6]](       Match[BLOCK:3:blocks:[7]]     ),     Match[BLOCK:2:blocks:[9]](       Match[BLOCK:3:blocks:[10]](         Match[BLOCK:4:blocks:[10]]       )     ),     Match[FUNCTION:1:blocks:[g]](       Match[BLOCK:2:blocks:[14]],       Match[BLOCK:2:blocks:[15]](         Match[BLOCK:3:blocks:[16]]       ),       Match[BLOCK:2:blocks:[18]]     )   ) ) </pre> |

La tabla 12.1 muestra el fuente y el resultado de `$blocks[0]->str` después de la llamada `my @blocks = $SimpleTrans::blocks->m($t);`. Obsérvese que los nodos `Match` son referencias a los nodos actuales del árbol.

### Usando el Depurador y str Para Ver los AST

Veamos un ejemplo mas de sesión con el depurador cuando la entrada proveída es el programa en la tabla 12.1:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> eyapp Simple4 ; perl -wd usesimple4.pl
```

```

main::(usesimple4.pl:6): $Data::Dumper::Indent = 1;
DB<1> f Simple4.eyp #<-- establecemos el fichero de trabajo
DB<2> l 643 #<-- nos paramos antes de la lectura
643: $self->YYData->{INPUT} = $_;
DB<3> b 643 $_=~ m{char d0} # <-- nos detendremos sólo si se trata de este programa
DB<4> c # <-- a correr!
..... # salida correspondiente a los fuentes anteriores
Simple4::Run(Simple4.eyp:643): $self->YYData->{INPUT} = $_;
DB<5> b Parse::Eyapp::Node::str # <-- Especificación completa de la subrutina en la que parar
DB<6> c
.....
Parse::Eyapp::Node::str(/home/pl/src/perl/YappWithDefaultAction/lib//Parse/Eyapp/Node.pm:543):
543: my $self = CORE::shift; # root of the subtree
DB<7> L # <-- Listar todos los breakpoints
/home/pl/src/perl/YappWithDefaultAction/lib//Parse/Eyapp/Node.pm:
543: my $self = CORE::shift; # root of the subtree
 break if (1)
Simple4.eyp:
643: $self->YYData->{INPUT} = $_;
 break if ($_=~ m{char d0})
DB<8> B * # <-- Eliminar todos los breakpoints
Deleting all breakpoints...
DB<9> b 654 # <-- detengámonos después del matching árbol
DB<10> c
Simple4::Run(Simple4.eyp:654): $_->node->{fatherblock} = $_->father->{node} for (@blocks
DB<4> x $Parse::Eyapp::Node::INDENT # <-- La variable INDENT controla el estilo en str
0 1
DB<11> x $t->str # <--- Salida con sangrado
0 '
PROGRAM(
 FUNCTION[f](
 BLOCK[3](
 BLOCK[4]
),
 BLOCK[6](
 BLOCK[7]
),
 BLOCK[9](
 BLOCK[10](
 BLOCK[10]
)
)
),
 FUNCTION[g](
 BLOCK[14],
 BLOCK[15](
 BLOCK[16]
),
 BLOCK[18]
)
),
DB<12> $Parse::Eyapp::Node::INDENT = 0 # <-- Estilo compacto. Por defecto
DB<13> x $t->str # <-- Todo en una línea

```

```

0 'PROGRAM(FUNCTION[f] (BLOCK[3] (BLOCK[4]),BLOCK[6] (BLOCK[7]),BLOCK[9] (BLOCK[10] (BLOCK[10]))),
FUNCTION[g] (BLOCK[14],BLOCK[15] (BLOCK[16]),BLOCK[18]))'
DB<14> . # Donde estoy? Que línea se va a ejecutar?
Simple4::Run(Simple4.eyp:654): $->node->{fatherblock} = $->father->{node} for (@blocks[1..$#
DB<15> c 655 # <-- Ejecutemos hasta la siguiente
Simple4::Run(Simple4.eyp:655): $Data::Dumper::Deepcopy = 1;
DB<16> x $blocks[0]->str # <-- Veamos la forma del nodo de matching
0 'Match[PROGRAM:0:blocks] (Match[FUNCTION:1:blocks:[f]] (Match[BLOCK:2:
blocks:[3]] (Match[BLOCK:3:blocks:[4]]),Match[BLOCK:2:blocks:[6]] (Match[
BLOCK:3:blocks:[7]]),Match[BLOCK:2:blocks:[9]] (Match[BLOCK:3:blocks:[1
0]] (Match[BLOCK:4:blocks:[10]]))),Match[FUNCTION:1:blocks:[g]] (Match[B
LOCK:2:blocks:[14]],Match[BLOCK:2:blocks:[15]] (Match[BLOCK:3:blocks:[1
6]]),Match[BLOCK:2:blocks:[18]]))' # <-- Todo está en una línea
DB<17> $Parse::Eyapp::Node::INDENT = 1
DB<18> x $blocks[0]->str
0 '
Match[PROGRAM:0:blocks] (
 Match[FUNCTION:1:blocks:[f]] (
 Match[BLOCK:2:blocks:[3]] (
 Match[BLOCK:3:blocks:[4]]
),
 Match[BLOCK:2:blocks:[6]] (
 Match[BLOCK:3:blocks:[7]]
),
 Match[BLOCK:2:blocks:[9]] (
 Match[BLOCK:3:blocks:[10]] (
 Match[BLOCK:4:blocks:[10]]
)
)
),
 Match[FUNCTION:1:blocks:[g]] (
 Match[BLOCK:2:blocks:[14]],
 Match[BLOCK:2:blocks:[15]] (
 Match[BLOCK:3:blocks:[16]]
),
 Match[BLOCK:2:blocks:[18]]
)
),
)'

```

## Variables que Gobiernan la Conducta de str

Las variables de paquete que gobiernan la conducta de **str** y sus valores por defecto aparecen en la siguiente lista:

- **our @PREFIXES = qw(Parse::Eyapp::Node::)**

La variable de paquete `Parse::Eyapp::Node::PREFIXES` contiene la lista de prefijos de los nombres de tipo que serán eliminados cuando **str** muestra el término. Por defecto contiene la cadena `'Parse::Eyapp::Node::'`. Es por esta razón que los nodos de matching que pertenecen a la clase `'Parse::Eyapp::Node::Match'` son abreviados a la cadena `Match` (línea 9).

- **our \$INDENT = 0** La variable `$Parse::Eyapp::Node::INDENT` controla el formato de presentación usado por `Parse::Eyapp::Node::str` :
  1. Si es 0 la representación es compacta.
  2. Si es 1 se usará un sangrado razonable.

3. Si es 2 cada paréntesis cerrar que este a una distancia mayor de `$Parse::Eyapp::Node::LINESEP` líneas será comentado con el tipo del nodo. Por defecto la variable `$Parse::Eyapp::Node::LINESEP` está a 4. Véase la página 699 para un ejemplo con `$INDENT` a 2 y `$Parse::Eyapp::Node::LINESEP` a 4.

- `our $STRSEP = ','`

Separador de nodos en un término. Por defecto es la coma.

- `our $DELIMITER = '['`

Delimitador de presentación de la información proveída por el método `info`. Por defecto los delimitadores son corchetes. Puede elegirse uno cualquiera de la lista:

`'[', '{', '(', '<'`

si se define como la cadena vacía o `undef` no se usarán delimitadores.

- `our $FOOTNOTE_HEADER = "\n-----\n"`

Delimitador para las notas a pie de página.

- `our $FOOTNOTE_SEP = ")\n"` Separador de la etiqueta/número de la nota al pie

- `our $FOOTNOTE_LEFT = '^{'`, `our $FOOTNOTE_RIGHT = '}'`

Definen el texto que rodea al número de referencia en el nodo.

## Footnotes

El siguiente ejemplo ilustra el manejo de *notas a pie de árbol* usando `str`. Se han definido los siguientes métodos:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code> sed -ne '677,$p' Simple6.eyp
```

```
$Parse::Eyapp::Node::INDENT = 1;
```

```
sub TERMINAL::info {
 my @a = join ':', @{$_[0]->{attr}};
 return "@a"
}
```

```
sub PROGRAM::footnote {
 return "Types:\n"
 .Dumper($_[0]->{types}).
 "Symbol Table:\n"
 .Dumper($_[0]->{symboltable})
}
```

```
sub FUNCTION::info {
 return $_[0]->{function_name}[0]
}
```

```
sub FUNCTION::footnote {
 return Dumper($_[0]->{symboltable})
}
```

```
sub BLOCK::info {
 return $_[0]->{line}
}
```

```

sub VAR::info {
 return $_[0]->{definition}{type} if defined $_[0]->{definition}{type};
 return "No declarado!";
}

```

```
*FUNCTIONCALL::info = *VARARRAY::info = \&VAR::info;
```

El resultado para el programa de entrada:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code> cat salida
int a,b;

```

```

int f(char c) {
 a[2] = 4;
 b[1][3] = a + b;
 c = c[5] * 2;
 return g(c);
}

```

produce una salida por stderr:

```
Identifier g not declared
```

y la siguiente descripción de la estructura de datos:

```

PROGRAM^{0}(
 FUNCTION{f}^{1}(
 ASSIGN(
 VARARRAY{INT}(
 TERMINAL{a:4},
 INDEXSPEC(
 INUM(
 TERMINAL{2:4}
)
)
),
 INUM(
 TERMINAL{4:4}
)
),
 ASSIGN(
 VARARRAY{INT}(
 TERMINAL{b:5},
 INDEXSPEC(
 INUM(
 TERMINAL{1:5}
),
 INUM(
 TERMINAL{3:5}
)
)
),
 PLUS(
 VAR{INT}(
 TERMINAL{a:5}
),

```





```

'CHAR' => bless({
 'children' => []
}, 'CHAR'),
'INT' => bless({
 'children' => []
}, 'INT')
};
Symbol Table:
$VAR1 = {
 'a' => {
 'type' => 'INT',
 'line' => 1
 },
 'b' => {
 'type' => 'INT',
 'line' => 1
 },
 'f' => {
 'type' => 'F(X_1(Char),INT)',
 'line' => 3
 }
};

```

```

1)
$VAR1 = {
 'c' => {
 'type' => 'CHAR',
 'param' => 1,
 'line' => 3
 }
};

```

## Usando str sobre una Lista de Árboles

Use `str` como método de clase

```
Parse::Eyapp::Node->str(@forest)
```

cuando quiera convertir a su representación término mas de un árbol. El código:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code> sed -ne '652,658p' Simple4.eypp \
| cat -n
1 local $Parse::Eyapp::Node::INDENT = 2;
2 local $Parse::Eyapp::Node::DELIMITER = "";
3 print $t->str."\n";
4 {
5 local $" = "\n";
6 print Parse::Eyapp::Node->str(@blocks)."\n";
7 }

```

produce la salida:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code> eyapp Simple4 ;\
treereg SimpleTrans.trg ;\

```

```

f() {
 int a,b[1][2],c[1][2][3];
 char d[10];
 b[0][1] = a;
}

PROGRAM(
 FUNCTION[f] (
 ASSIGN(
 VARARRAY(
 TERMINAL[b:4],
 INDEXSPEC(
 INUM(
 TERMINAL[0:4]
),
 INUM(
 TERMINAL[1:4]
)
) # INDEXSPEC
) # VARARRAY,
 VAR(
 TERMINAL[a:4]
)
) # ASSIGN
) # FUNCTION
) # PROGRAM

Match[PROGRAM:0:blocks] (
 Match[FUNCTION:1:blocks:[f]]
)

```

Obsérvense los comentarios # TIPODENODO que acompañan a los paréntesis cerrar cuando estos están lejos (esto es, a una distancia de mas de \$Parse::Eyapp::Node::LINESEP líneas) del correspondiente paréntesis abrir. Tales comentarios son consecuencia de haber establecido el valor de \$Parse::Eyapp::Node::INDENT a 2.

## 12.9. Práctica: Establecimiento de la relación uso-declaración

En la segunda parte del análisis de ámbito se resuelve el problema de la identificación de los nombres estableciendo la relación entre el uso de un nombre y la declaración que se le aplica (dada por los atributos **scope** y **type**):

Añada las siguientes comprobaciones dependientes del contexto al compilador de Simple C:

- Resuelva el problema del ámbito de los identificadores: asocie cada uso de un identificador con la declaración que se le aplica.
- Asocie cada nodo RETURN con el nodo de la subrutina que lo engloba. Añádale un atributo al nodo RETURN que indique el tipo que la rutina debe retornar.
- Asocie cada nodo CONTINUE y cada nodo BREAK con el nodo del bucle que lo engloba. Indique el error si la sentencia (CONTINUE o BREAK) no esta dentro de un bucle. Añada a dichos nodos un atributo que referencia al nodo del bucle que lo engloba.

- Añada etiquetas de salto y la sentencia `goto` al lenguaje Simple C. Toda etiqueta va asociada con una sentencia. La sintaxis de statement debe ser extendida para incluir sentencias etiquetadas:

```

statement:
 expression ';' { $_[1] }
| ';'
| %name BREAK
 $BREAK ';'
 { ... }
| %name CONTINUE
 $CONTINUE ';'
 { ... }
| 'goto' ID ';'
| ID ':' statement
 { ... }
;

```

Las siguientes indicaciones tienen por objeto ayudarle en el análisis de ámbito de etiquetas:

- Sólo se permiten saltos dentro de la misma subrutina.
- Será necesario crear un manejador de ámbito de etiquetas mediante `Parse::Eyapp::Scope->new`
- Al comienzo de cada función comienza un ámbito de etiquetas.
- Cada vez que se encuentra una sentencia etiquetada la etiqueta ID es guardada en una tabla de símbolos de etiquetas.
- El valor asociado con la entrada para la etiqueta es una referencia a un hash. La clave mas importante en ese hash es la referencia a la sentencia que etiqueta.
- Cada aparición de un goto es una instancia de ámbito.
- La llamada a `end_scope` (al terminar el análisis de la función) deberá decorar automáticamente el nodo GOTO con la referencia a la sentencia asociada con la etiqueta
- La tabla de etiquetas es un atributo del nodo asociado con la función.

## 12.10. Práctica: Establecimiento de la Relación Uso-Declaración Usando Expresiones Regulares Árbol

En la práctica anterior resolvió los problemas de ámbito usando `Parse::Eyapp::Scope`. En esta se pide resolverlos usando expresiones regulares árbol:

- Asocie cada nodo `RETURN` con el nodo de la subrutina que lo engloba. Añádale un atributo al nodo `RETURN` que indique el tipo que la rutina debe retornar.
- Asocie cada nodo `CONTINUE` y cada nodo `BREAK` con el nodo del bucle que lo engloba. Indique el error si la sentencia (`CONTINUE` o `BREAK`) no esta dentro de un bucle. Añada a dichos nodos un atributo que referencia al nodo del bucle que lo engloba.
- Añada etiquetas de salto y la sentencia `goto` al lenguaje Simple C. Toda etiqueta va asociada con una sentencia.

## 12.11. Práctica: Estructuras y Análisis de Ámbito

Los registros o estructuras introducen espacios de nombres. Los selectores de campo viven en los espacios de nombres asociados con los registros. Modifique de manera apropiada la gramática de `SimpleC` para introducir el tipo de datos `struct`. Complete el análisis sintáctico.

- Es necesario introducir nuevas variables sintácticas:

```

type_specifier:
 basictype
 | struct
;

struct:
 'struct' ID '{' struct_decl_list '}'
 | 'struct' ID
;

struct_decl_list: (type_specifier declList ';') +
;

```

Recuerde que:

```

declList: (ID arraySpec) <%+ ','>
;

```

- Es necesario también modificar algunas categorías gramaticales existentes, como **declaration**:

```

declaration:
 %name DECLARATION
 $type_specifier $declList ';'
;

```

También es necesario modificar **Variable**:

```

Variable:
 ID
 | Variable ('[' binary ']')
 | Variable '.' ID
;

```

Sigue un ejemplo de derivación generando una variable de la forma **a[i].x[4].b**:

```

Variable => Variable.ID => Variable[binary].ID
=> Variable.ID[binary].ID => Variable[binary].ID[binary].ID
=> ID[binary].ID[binary].ID

```

Veamos otro ejemplo de derivación generando una variable de la forma **a[i].x[4]**:

```

Variable => Variable[binary] => Variable.ID[binary] => Variable[binary].ID[binary]
=> ID[binary].ID[binary]

```

- Es necesario extender las expresiones de tipo para incluir las **struct**. Un registro es similar al producto cartesiano (usado en las expresiones de tipo de las funciones), sólo que los campos tienen nombres. Así:

```

struct row {
 int address;
 char lexeme[10];
} table[100];

```

la expresión de tipo para `struct row` podría ser:

```
STRUCT_row(address, INT, lexeme, A_10(CHAR))
```

o bien:

```
STRUCT_row(X_2(address, INT), X_2(lexeme, A_10(CHAR)))
```

Como se ve en el ejemplo el identificador de campo (`address, lexeme`) forma parte de la descripción del tipo.

- La introducción de las `struct` produce una cierta sobrecarga de identificadores. El siguiente programa es correcto:

```
1 struct x {
2 int z;
3 int b;
4 } x;
5
6 main() {
7 x;
8 }
```

Aquí `x` es usado como nombre de un tipo y como nombre de variable. Plantee soluciones al problema del almacenamiento en la tabla de símbolos. Una posibilidad es que la entrada en la tabla de símbolos para el tipo sea `struct x` mientras que para la variable es `x`.

- El cálculo de que declaración se aplica a una `struct` puede hacerse de manera similar a como ha sido ilustrado en las prácticas anteriores. Otra cosa son los identificadores de campo. Por ejemplo en

```
a[i+j].x = 4;
```

Es necesario calcular que `a[i+j]` tiene tipo `struct` y comprobar en ese momento que dicha `struct` tiene un campo `x` que es de tipo `int`. Hasta que no hemos calculado el tipo asociado con la expresión `a[i+j]` no sabemos si el uso de `x` es legal o no. Realmente el ámbito del identificador `x` no se hace en tiempo de análisis de ámbito sino en tiempo de comprobación de tipos. El cálculo del ámbito para estos identificadores de campo se mezcla con el cálculo de tipos. En cierto modo el uso de `x` en `a[i+j].x = 4` forma parte de la descripción del tipo:

```
STRUCT_point(X_2(x, INT), X_2(y, INT))
```

Por tanto en las reglas de `Variable`:

`Variable`:

```
 ID
 | Variable ('[' binary ']')
 | Variable '.' ID
;
```

A estas alturas solo estamos en condiciones de asociarle su definición al identificador principal, esto es, al de la regla  $Variable \rightarrow ID$ . Los demás deberán esperar a la fase de análisis de tipos.

- Hasta ahora casi lo único que guardabamos en las tablas de símbolos era el tipo y el número de línea ( otros valores posibles son `param` para determinar si es un parámetro, o una referencia a la tabla de símbolos de la funcion en el caso de las funciones). Esto era suficiente porque sólo teníamos dos tipos de identificadores: los de variable y los de función. Ahora que tenemos tres tipos: funciones, variables y tipos es conveniente añadir un atributo `kind` que guarde la clase de objeto que es.
- C requiere que los nombres de tipo sean declarados antes de ser usados (con la excepción de las referencias a tipos `struct` no declarados). Por ejemplo, el siguiente programa:

```

1 struct x {
2 struct y z;
3 int b;
4 } w;
5
6 struct y {
7 int a;
8 };

```

produce un mensaje de error:

```
prueba1.c:2: error: field 'z' has incomplete type
```

Esto impide definiciones recursivas como:

```

1 struct x {
2 struct y z;
3 int b;
4 } w;
5
6 struct y {
7 struct x a;
8 };

```

A estas alturas del análisis estamos en condiciones de comprobar que no existan usos de declaraciones no definidas.

- Para simplificar el análisis de tipos posterior, las funciones en Simple C devuelven un tipo básico. En C las funciones pueden devolver un `struct` aunque no pueden devolver un array ni una función. Si que pueden sin embargo devolver un puntero a un array o a una función.

La restricción que imponemos de que el tipo retornado sea un tipo básico no significa que necesariamente la regla sintáctica deba especificar tal condición.

- Para simplificar la práctica consideraremos que todos los identificadores de tipos `struct` son globales. Por tanto es un error tener dos declaraciones de tipo `struct` con el mismo nombre. Claro está, si desea implantar ámbitos en los identificadores de tipos puede hacerlo.
- Consulte las siguientes fuentes:

- Para ver una gramática de C consulte la página [Cgrammar.html](http://www.lysator.liu.se/c/ANSI-C-grammar-y.html).
- Una gramática de ANSI C se encuentra en <http://www.lysator.liu.se/c/ANSI-C-grammar-y.html>.

## Capítulo 13

# Análisis de Tipos

### 13.1. Análisis de Tipos: Conceptos Básicos

En la mayoría de los lenguajes los objetos manipulados son declarados en alguna parte del programa y usados en otras. Ya dijimos que el análisis de ámbito es el cálculo de la función que asigna a un uso de un objeto la definición que se le aplica.

El análisis de tipos tiene por objetivo asegurar que el uso de los objetos definidos es correcto: esto es, que su uso se atiene a la semántica de su definición; por ejemplo, que un array de enteros no es llamado como función o que no se intenta incrementar una función o que el valor retornado por una función es de la naturaleza descrita en su definición.

#### Expresiones de Tipo, Sistemas de Tipos y Comprobadores de Tipos

**Definición 13.1.1.** *Una forma adecuada de representar los tipos dentro de un compilador es usando un lenguaje de expresiones de tipo. Un lenguaje de las expresiones de tipo debe describir de manera clara y sencilla los tipos del lenguaje fuente. No confunda este lenguaje con el sub-lenguaje del lenguaje fuente que consiste en las declaraciones o definiciones. No tienen por que ser iguales. El compilador traduce las declaraciones de tipo en expresiones de tipo. El lenguaje de las expresiones de tipo es la representación interna que el compilador tiene de estas declaraciones y depende del compilador. El lenguaje de las declaraciones no.*

**Definición 13.1.2.** *Un sistema de tipos de un lenguaje/compilador es el conjunto de reglas del lenguaje (que es traducido e interpretado por el compilador que permiten asignar expresiones de tipo a las instancias de uso de los objetos del programa.*

Si bien el sistema de tipos es una propiedad del lenguaje, no es raro que los compiladores introduzcan modificaciones en el sistema de tipos del lenguaje. Por ejemplo en Pascal el tipo de un array incluye los índices del array<sup>1</sup>). Esto y las reglas de equivalencia de tipos de Pascal limitan gravemente la genericidad de las funciones en Pascal. Por eso algunos compiladores Pascal permiten en una llamada a función la compatibilidad de tipos entre arrays de diferente tamaño y diferentes conjuntos de índices. Desgraciadamente la forma en la que lo hacen puede diferir de compilador a compilador.

**Definición 13.1.3.** *Un comprobador de tipos verifica que el uso de los objetos en los constructos de uso se atiene a lo especificado en sus declaraciones o definiciones de acuerdo a las reglas especificadas por el sistema de tipos.*

#### Tipado Estático y Tipado Dinámico

**Definición 13.1.4.** *Un lenguaje de programación tiene tipado estático si su comprobación de tipos ocurre en tiempo de compilación sin tener que comprobar equivalencias en tiempo de ejecución.*

*Un lenguaje de programación tiene tipado dinámico si el lenguaje realiza comprobaciones de tipo en tiempo de ejecución. En un sistema de tipos dinámico los tipos suelen estar asociados con los valores no con las variables.*

---

<sup>1</sup>Traduttore, traditore!



**Definición 13.1.5.** *El tipado dinámico hace mas sencilla la escritura de metaprogramas: programas que reciben como datos otros códigos y los manipulan para producir nuevos códigos. `Parse::EYapp` y `Parse::Tree::Regexp` son ejemplos de metaprogramación. Cada vez que escribimos un procesador de patrones, templates o esqueletos de programación estamos haciendo meta-programación.*

*El lenguaje en el que se escribe el metaprograma se denomina metalenguaje. El lenguaje al que se traduce el metaprograma se denomina lenguaje objeto. La capacidad de un lenguaje de programación para ser su propio metalenguaje se denomina reflexividad. Para que haya reflexión es conveniente que el código sea un tipo de estructura de datos soportado por el lenguaje al mismo nivel que otros tipos básicos y que sea posible traducir dinámicamente texto a código.*

## Tipado Fuerte y Tipado Débil

**Definición 13.1.6.** *Aunque el significado de los términos Fuertemente Tipado y su contrario Débilmente Tipado varían con los autores, parece haber consenso en que los lenguajes con tipado fuerte suelen reunir alguna de estas características:*

- *La comprobación en tiempo de compilación de las violaciones de las restricciones impuestas por el sistema de tipos. El compilador asegura que para cualesquiera operaciones los operandos tienen los tipos válidos.*
- *Toda operación sobre tipos inválidos es rechazada bien en tiempo de compilación o de ejecución.*
- *Algunos autores consideran que el término implica desactivar cualquier conversión de tipos implícita. Si el programador quiere una conversión deberá explicitarla.*
- *La ausencia de modos de evadir al sistema de tipos.*
- *Que el tipo de un objeto de datos no varíe durante la vida del objeto. Por ejemplo, una instancia de una clase no puede ver su clase alterada durante la ejecución.*

## Sobrecarga, Polimorfismo e Inferencia de Tipos

Un símbolo se dice *sobrecargado* si su significado varía dependiendo del contexto. En la mayoría de los lenguajes Los operadores aritméticos suelen estar sobrecargados, dado que se sustancian en diferentes algoritmos según sus operandos sean enteros, flotantes, etc.

En algunos lenguajes se permite la sobrecarga de funciones. así es posible tener dos funciones llamadas `min`:

```
int min(int a, int b) {
 if (a < b) return a;
 return b;
}
```

```
int min(string a, string b) {
 if (strcmp(a, b) < 0) return a;
 return b;
}
```

A la hora de evaluar el tipo de las expresiones es el contexto de la llamada el que determina el tipo de la expresión:

```
float x,y;
int a,b;
string c,d;
```

```
u = min(x,y); /* Puede que correcto: x e y seran truncados a enteros. Tipo entero */
v = min(a,b); /* Correcto: Tipo devuelto es entero */
w = min(c,d); /* Correcto: Tipo devuelto es string */
t = min(x,c); /* Error */
```

### Ejercicio 13.1.1. ¿Como afecta al análisis de ámbito la sobrecarga de operadores?

**Definición 13.1.7.** La inferencia de tipos hace referencia a aquellos algoritmos que deducen automáticamente en tiempo de compilación - sin información adicional del programador, o bien con anotaciones parciales del programador - el tipo asociado con un uso de un objeto del programa. Un buen número de lenguajes de programación funcional permiten implantar inferencia de tipos (Haskell, OCaml, ML, etc).

Véase como ejemplo de inferencia de tipos la siguiente sesión en OCaml:

```
pl@nereida:~/src/perl/attributegrammar/Language-AttributeGrammar-0.08/examples$ ocaml
Objective Caml version 3.09.2
```

```
let minimo = fun i j -> if i<j then i else j;;
val minimo : 'a -> 'a -> 'a = <fun>
minimo 2 3;;
- : int = 2
minimo 4.9 5.3;;
- : float = 4.9
minimo "hola" "mundo";;
- : string = "hola"
```

El compilador OCaml infiere el tipo de las expresiones. Así el tipo asociado con la función `minimo` es `'a -> 'a -> 'a` que es una *expresión de tipo* que contiene *variables de tipo*. El operador `->` es asociativo a derechas y así la expresión debe ser leída como `'a -> ('a -> 'a)`. Básicamente dice: es una función que toma un argumento de tipo `'a` (donde `'a` es una variable tipo que será instanciada en el momento del uso de la función) y devuelve una función que toma elementos de tipo `'a` y retorna elementos de tipo `'a`.

**Definición 13.1.8.** Aunque podría pensarse que una descripción mas adecuada del tipo de la función `minimo` fuera `'a x 'a -> 'a`, lo cierto es que en algunos lenguajes funcionales es usual que todas las funciones sean consideradas como funciones de una sola variable. La función de dos variables `'a x 'a -> 'a` puede verse como una función `'a -> ('a -> 'a)`. En efecto la función `minimo` cuando recibe un argumento retorna una función:

```
let min_mundo = minimo "mundo";;
val min_mundo : string -> string = <fun>
min_mundo "pedro";;
- : string = "mundo"
min_mundo "antonio";;
- : string = "antonio"
min_mundo 4;;
This expression has type int but is here used with type string
min_mundo(string_of_int(4));;
- : string = "4"
```

Esta estrategia de reducir funciones de varias variables a funciones de una variable que retornan funciones de una variable se conoce con el nombre de *currying* o aplicación parcial.

**Definición 13.1.9.** El polimorfismo es una propiedad de ciertos lenguajes que permite una interfaz uniforme a diferentes tipos de datos.

Se conoce como *función polimorfa* a una función que puede ser aplicada o evaluada sobre diferentes tipos de datos.

Un tipo de datos se dice *polimorfo* si es un tipo de datos generalizado o no completamente especificado. Por ejemplo, una lista cuyos elementos son de cualquier tipo.

**Definición 13.1.10.** Se llama Polimorfismo Ad-hoc a aquel en el que el número de combinaciones que pueden usarse es finito y las combinaciones deben ser definidas antes de su uso. Se habla de polimorfismo paramétrico si es posible escribir el código sin mención específica de los tipos, de manera que el código puede ser usado con un número arbitrario de tipos.

Por ejemplo, la herencia y la sobrecarga de funciones y métodos son mecanismos que proveen polimorfismo ad-hoc. Los lenguajes funcionales, como OCaml suelen proveer polimorfismo paramétrico. En OOP el polimorfismo paramétrico suele denominarse *programación genérica*

En el siguiente ejemplo en OCaml construimos una función similar al `map` de Perl. La función `mymap` ilustra el polimorfismo paramétrico: la función puede ser usada con un número arbitrario de tipos, no hemos tenido que hacer ningún tipo de declaración explícita y sin embargo el uso incorrecto de los tipos es señalado como un error:

```
let rec mymap f list =
 match list with
 [] -> []
 | hd :: tail -> f hd :: mymap f tail;;
val mymap : ('a -> 'b) -> 'a list -> 'b list = <fun>
mymap (function n -> n*2) [1;3;5];;
- : int list = [2; 6; 10]
mymap (function n -> n.[0]) ["hola"; "mundo"];;
- : char list = ['h'; 'm']
mymap (function n -> n*2) ["hola"; "mundo"];;
This expression has type string but is here used with type int
```

### Equivalencia de Expresiones de Tipo

La introducción de nombres para las expresiones de tipo introduce una ambigüedad en la interpretación de la equivalencia de tipos. Por ejemplo, dado el código:

```
typedef int v10[10];
v10 a;
int b[10];
```

¿Se considera que `a` y `b` tienen tipos compatibles?

**Definición 13.1.11.** Se habla de equivalencia de tipos estructural cuando los nombres de tipo son sustituidos por sus definiciones y la equivalencia de las expresiones de tipo se traduce en la equivalencia de sus árboles sintácticos o DAGs. Si los nombres no son sustituidos se habla de equivalencia por nombres o de equivalencia de tipos nominal.

Si utilizamos la opción de sustituir los nombres por sus definiciones y permitimos en la definición de tipo el uso de nombres de tipo no declarados se pueden producir ciclos en el grafo de tipos.

El lenguaje C impide la presencia de ciclos en el grafo de tipos usando dos reglas:

1. Todos los identificadores de tipo han de estar definidos antes de su uso, con la excepción de los punteros a registros no declarados
2. Se usa equivalencia estructural para todos los tipos con la excepción de las `struct` para las cuales se usa equivalencia nominal

Por ejemplo, el siguiente programa:

```
nereida:~/src/perl/testing> cat -n typeequiv.c
1 #include <stdio.h>
2
3 typedef struct {
4 int x, y;
```

```

5 struct record *next;
6 } record;
7
8 record z,w;
9
10 struct recordcopy {
11 int x, y;
12 struct recordcopy *next;
13 } r,k;
14
15
16 main() {
17 k = r; /* no produce error */
18 z = w; /* no produce error */
19 r = z;
20 }

```

Produce el siguiente mensaje de error:

```

nereida:~/src/perl/testing> gcc -fsyntax-only typeequiv.c
typeequiv.c: En la función 'main':
typeequiv.c:19: error: tipos incompatibles en la asignación

```

En lenguajes dinámicos una forma habitual de equivalencia de tipos es el tipado pato:

**Definición 13.1.12.** *Se denomina duck typing o tipado pato a una forma de tipado dinámico en la que el conjunto de métodos y propiedades del objeto determinan la validez de su uso. Esto es: dos objetos pertenecen al mismo tipo-pato si implementan/soportan la misma interfaz independientemente de si tienen o no una relación en la jerarquía de herencia.*

*El término hace referencia al llamado test del pato: If it waddles like a duck, and quacks like a duck, it's a duck!.*

## 13.2. Conversión de Tipos

El comprobador de tipos modifica el árbol sintáctico para introducir *conversiones de tipo* donde sean necesarias. Por ejemplo, si tenemos

```

int i;
float x;

x+i;

```

Dado el árbol de la expresión PLUS(VAR, VAR), el analizador de tipos introducirá un nodo intermedio INT2FLOAT para indicar la necesidad de la conversión y especificará el tipo de PLUS que se usa:

PLUSFLOAT(VAR, INT2FLOAT(VAR)).

Una transformación árbol de optimización que entra en este punto es la conversión de tipo en tiempo de compilación de las constantes. Por ejemplo, dados los dos programas:

|                                                                          |                                                                            |
|--------------------------------------------------------------------------|----------------------------------------------------------------------------|
| <pre> float X[N]; int i;  for(i=0; i&lt;N; i++) {     X[i] = 1; } </pre> | <pre> float X[N]; int i;  for(i=0; i&lt;N; i++) {     X[i] = 1.0; } </pre> |
|--------------------------------------------------------------------------|----------------------------------------------------------------------------|

los efectos sobre el rendimiento serán lamentables si el compilador no realiza la conversión de la constante entera 1 del programa de la izquierda en tiempo de compilación sino que la conversión se deja a una subrutina de conversión que es llamada en tiempo de ejecución. En tal caso se obtendrían rendimientos completamente diferentes para los programas en la izquierda y en la derecha.

**Ejercicio 13.2.1.** *Escriba el pseudocódigo de las transformaciones árbol Treeregexp que implementan la conversión de tipos de constantes en tiempo de compilación*

### 13.3. Expresiones de Tipo en Simple C

En las secciones anteriores dedicadas al análisis de ámbito introdujimos nuestro lenguaje de expresiones de tipo.

Así la expresión `F(X_2(INT,CHAR),INT)` describe el tipo de las funciones de 2 parámetros que devuelven enteros. El primer parámetro debe ser entero y el segundo un carácter.

Otro ejemplo es la expresión de tipo `A_10(A_20(INT))` que describe un vector 10x20 de enteros.

**Ejercicio 13.3.1.** *Describe la gramática independiente del contexto que genera las expresiones de tipo usadas por el compilador de Simple C tal y como fueron presentadas en las secciones previas*

Observe que en tanto en cuanto las expresiones de tipo son frases de un lenguaje - cuya semántica es describir los tipos del lenguaje fuente - tiene un árbol sintáctico asociado. En las secciones anteriores sobre el análisis de ámbito presentamos una versión inicial de las expresiones de tipo y las utilizamos mediante una representación dual: como cadena y como árbol. Para ello usábamos el método `Parse::Eyapp::Node->new` el cual nos permite obtener una representación árbol de la expresión de tipo. En un contexto escalar este método devuelve la referencia al árbol sintáctico construido:

```
nereida:~/src/perl/testing> perl -MParse::Eyapp::Node -de 0
Loading DB routines from perl5db.pl version 1.28
Editor support available.
main::(-e:1): 0
DB<1> x scalar(Parse::Eyapp::Node->new('F(X_2(INT,CHAR),INT)'))
0 F=HASH(0x859fb58)
 'children' => ARRAY(0x85d52ec)
 0 X_2=HASH(0x859fb70)
 'children' => ARRAY(0x859fb28)
 0 INT=HASH(0x85a569c)
 'children' => ARRAY(0x852cc84)
 empty array
 1 CHAR=HASH(0x859fb34)
 'children' => ARRAY(0x850a7cc)
 empty array
 1 INT=HASH(0x859fab0)
 'children' => ARRAY(0x852c9d8)
 empty array
```

### 13.4. Construcción de las Declaraciones de Tipo con hnew

Una manera aún mas conveniente de representar una expresión de tipo es a través de un *grafo dirigido acíclico* o *DAG* (*Directed Acyclic Graph*).

`hnew`

`Parse::Eyapp` provee el método `Parse::Eyapp::Node->hnew` el cual facilita la construcción de DAGs. El método `hnew` trabaja siguiendo lo que en programación funcional se denomina *hashed consing*. En *Lisp* la función de asignación de memoria es `cons` y la forma habitual de compartición es

a través de una tabla hash. Hash consing es una técnica para compartir datos que son estructuralmente iguales. De esta forma se ahorra memoria. La técnica esta relacionada con - puede considerarse una aplicación concreta de - la técnica conocida con el nombre de *memoization* [?]. Si alguno de los subárboles (o el árbol) fueron creados anteriormente no serán creados de nuevo sino que `hnew` retorna una referencia al ya existente. Esto hace mas rápida la comprobación de tipos: *en vez de comprobar que los dos árboles son estructuralmente equivalentes basta con comprobar que las referencias son iguales*. Así la comprobación de la equivalencia de tipos pasa a requerir tiempo constante.

Observe los siguientes comandos:

```
DB<2> x $a = Parse::Eyapp::Node->hnew('F(X_3(A_3(A_5(INT)), CHAR, A_5(INT)),CHAR)')
0 F=HASH(0x85f6a20)
 'children' => ARRAY(0x85e92e4)
 |- 0 X_3=HASH(0x83f55fc)
 | 'children' => ARRAY(0x83f5608)
 | |- 0 A_3=HASH(0x85a0488)
 | | 'children' => ARRAY(0x859fad4)
 | | 0 A_5=HASH(0x85e5d3c)
 | | 'children' => ARRAY(0x83f4120)
 | | 0 INT=HASH(0x83f5200)
 | | 'children' => ARRAY(0x852ccb4)
 | | empty array
 | |- 1 CHAR=HASH(0x8513564)
 | | 'children' => ARRAY(0x852cad4)
 | | empty array
 | '- 2 A_5=HASH(0x85e5d3c)
 | -> REUSED_ADDRESS
 '- 1 CHAR=HASH(0x8513564)
 -> REUSED_ADDRESS
DB<3> x $a->str
0 'F(X_3(A_3(A_5(INT)),CHAR,A_5(INT)),CHAR)'
```

La segunda aparición de `A_5(INT)` aparece etiquetada como reusada (tercer hijo de `X_3`). Lo mismo ocurre con la segunda aparición de `CHAR`.

### **hnew y el Manejador de Atributos**

El valor asignado a los atributos no influye en la memoización de `hnew`: si los términos tienen la misma forma se devolverán los mismos árboles:

```
pl@nereida:~/src/perl/testing$ perl -MParse::Eyapp::Node -wde 0
main::(-e:1): 0
DB<1> @x = Parse::Eyapp::Node->hnew('A(B,C(D))', sub { $_->{n} = $i++ for @_ })
DB<2> @y = Parse::Eyapp::Node->hnew('C(D)', sub { $_->{t} = 10-$j++ for @_ })
DB<3> x @x
0 A=HASH(0x850935c)
 'children' => ARRAY(0x8509530)
 0 B=HASH(0x85095a8)
 'children' => ARRAY(0x8326800)
 empty array
 'n' => 1
 1 C=HASH(0x8509410)
 'children' => ARRAY(0x85094dc)
 0 D=HASH(0x85094b8)
 'children' => ARRAY(0x85095f0)
 empty array
```

```

 'n' => 3
 't' => 9
 'n' => 2
 't' => 10
 'n' => 0
1 B=HASH(0x85095a8)
 -> REUSED_ADDRESS
2 C=HASH(0x8509410)
 -> REUSED_ADDRESS
3 D=HASH(0x85094b8)
 -> REUSED_ADDRESS
DB<4> x @y
0 C=HASH(0x8509410)
 'children' => ARRAY(0x85094dc)
 0 D=HASH(0x85094b8)
 'children' => ARRAY(0x85095f0)
 empty array
 'n' => 3
 't' => 9
 'n' => 2
 't' => 10
1 D=HASH(0x85094b8)
 -> REUSED_ADDRESS

```

### La Función build\_function\_scope

Sigue el código de la función build\_function\_scope:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/sub bu.*scop/,/~/p' Types.eypp | cat -n
1 sub build_function_scope {
2 my ($funcDef, $returntype) = @_ ;
3
4 my $function_name = $funcDef->{function_name}[0];
5 my @parameters = @{$funcDef->{parameters}};
6 my $lst = $funcDef->{symboltable};
7 my $numargs = scalar(@parameters);
8
9 #compute type
10 my $partype = "";
11 my @types = map { $lst->{$_}{type} } @parameters;
12 $partype .= join ",", @types if @types;
13
14 my $type = "F(X_$numargs($partype),$returntype)";
15
16 #insert it in the hash of types
17 $type{$type} = Parse::Eyapp::Node->hnew($type);
18 $funcDef->{type} = $type;
19 $funcDef->{t} = $type{$type};
20
21 #insert it in the global symbol table
22 die "Duplicated declaration of $function_name at line $funcDef->{function_name}[1]\n"
23 if exists($st{$function_name});
24 $st{$function_name}->{type} = $type;
25 $st{$function_name}->{line} = $funcDef->{function_name}[1];

```

```

26
27 return $funcDef;
28 }

```

Obsérvese que para acelerar los accesos al atributo tipo dotamos a los nodos del atributo `t` que permite acceder directamente al DAG representando el tipo.

## 13.5. Inicio de los Tipos Básicos

Al comienzo y al final del análisis del programa fuente se inician mediante la llamada a `reset_file_scope_vars` los tipos básicos `INT`, `CHAR` y `VOID`. Este último se utilizará en las expresiones de tipo para representar la ausencia de tipo.

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/^program:\/,\/^;/p' Types.eypp | cat -n
1 program:
2 {
3 reset_file_scope_vars();
4 }
5 definition<%name PROGRAM +>.program
6 {
7 $program->{symboltable} = { %st }; # creates a copy of the s.t.
8 $program->{depth} = 0;
9 $program->{line} = 1;
10 $program->{types} = { %type };
11 $program->{lines} = $tokenend;
12
13 my ($nondec, $declared) = $ids->end_scope($program->{symboltable}, $program, 'type'
14
15 if (@$nondec) {
16 warn "Identifier ".$_->key." not declared at line ".$_->line."\n" for @$nondec;
17 die "\n";
18 }
19
20 # Type checking: add a direct pointer to the data-structure
21 # describing the type
22 $_->{t} = $type{$_->{type}} for @$declared;
23
24 my $out_of_loops = $loops->end_scope($program);
25 if (@$out_of_loops) {
26 warn "Error: ".ref($_)." outside of loop at line $_->{line}\n" for @$out_of_loop
27 die "\n";
28 }
29
30 # Check that are not dangling breaks
31 reset_file_scope_vars();
32
33 $program;
34 }
35 ;

```

Se procede a establecer el atributo `t` como enlace directo a la expresión de tipo. Lo mismo se hace en las funciones y en los bloques.

Sigue el código de `reset_file_scope_vars`:



```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/^sub re.*vars/,/^}/p' Types.eypp | cat -n
1 sub reset_file_scope_vars {
2 %st = (); # reset symbol table
3 ($tokenbegin, $tokenend) = (1, 1);
4 %type = (INT => Parse::Eyapp::Node->hnew('INT'),
5 CHAR => Parse::Eyapp::Node->hnew('CHAR'),
6 VOID => Parse::Eyapp::Node->hnew('VOID'),
7);
8 $depth = 0;
9 $ids = Parse::Eyapp::Scope->new(
10 SCOPE_NAME => 'block',
11 ENTRY_NAME => 'info',
12 SCOPE_DEPTH => 'depth',
13);
14 $loops = Parse::Eyapp::Scope->new(
15 SCOPE_NAME => 'exits',
16);
17 $ids->begin_scope();
18 $loops->begin_scope(); # just for checking
19 }

```

## 13.6. Comprobación Ascendente de los Tipos

La comprobación de tipos se hace mediante una visita ascendente del AST: primero se comprueban y computan los tipos de los hijos del nodo aplicándole las reglas del sistema de tipos que hemos implantado usando transformaciones árbol. Después pasamos a computar el tipo del nodo. Para la visita usamos el método `bud` (por bottom-up decorator<sup>2</sup>) de `Parse::Eyapp::Node`.

La llamada `$t->bud(@typecheck)` hace que se visite cada uno de los nodos de `$t` en orden bottom-up. Para cada nodo se busca una transformación árbol en la lista `@typecheck` que se pueda aplicar. Tan pronto como se encuentra una se aplica y se procede con el siguiente nodo.

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/^sub compile/,/^}/p' Types.eypp | cat -n
1 sub compile {
2 my($self)=shift;
3
4 my ($t);
5
6 $self->YYData->{INPUT} = $_[0];
7
8 $t = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
9 #yydebug => 0x1F
10);
11
12 # Scope Analysis: Block Hierarchy
13 our $blocks;
14 my @blocks = $blocks->m($t);
15 $_->node->{fatherblock} = $_->father->{node} for (@blocks[1..$#blocks]);

```

---

<sup>2</sup>De acuerdo al diccionario:

`bud`: n. Botany. A small protuberance on a stem or branch, sometimes enclosed in protective scales and containing an undeveloped shoot, leaf, or flower.

```

16
17 # Scope Analysis: Return-Function
18 our $retscope; # retscope: /FUNCTION|RETURN/
19 my @returns = $retscope->m($t);
20 for (@returns) {
21 my $node = $_->node;
22 if (ref($node) eq 'RETURN') {
23 my $function = $_->father->node;
24 $node->{function} = $function;
25 $node->{t} = $function->{t}->child(1);
26 }
27 }
28
29 # Type checking
30 set_types($t); # Init basic types
31
32 my @typecheck = (# Check typing transformations for
33 our $inum, # - Numerical constantss
34 our $charconstant, # - Character constants
35 our $bin, # - Binary Operations
36 our $arrays, # - Arrays
37 our $assign, # - Assignments
38 our $control, # - Flow control sentences
39 our $functioncall, # - Function calls
40 our $statements, # - Those nodes with void type
41 # (STATEMENTS, PROGRAM, etc.)
42 our $returntype, # - Return
43);
44
45 $t->bud(@typecheck);
46
47 # The type checking for trees RETURN exp is made
48 # in adifferent way. Just for fun
49 #our $bind_ret2function;
50 #my @FUNCTIONS = $bind_ret2function->m($t);
51
52 return $t;
53 }

```

La llamada a la subrutina `set_types($t)` tiene por objeto establecer la comunicación entre el programa árbol en `Trans.trg` y el compilador en `Types.eyp`. La subrutina se encuentra en el código de apoyo en el programa árbol:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '17,62p' Trans.trg | cat -n
1 {
2
3 my $types; # reference to the hash containing the type table
4 my ($INT, $CHAR, $VOID);
5
6 sub type_error {
7 my $msg = shift;
8 my $line = shift;
9 die "Type Error at line $line: $msg\n"

```

```

10 }
11
12 sub set_types {
13 my $root = shift;
14 $types = $root->{types};
15 $INT = $types->{INT};
16 $CHAR = $types->{CHAR};
17 $VOID = $types->{VOID};
18 }
19
20 sub char2int {
21 ;
22 }
23
24 sub int2char {
25 ;
26 }
27
28 }

```

La subrutina `set_types` inicializa la variable léxica `$types` que referencia a la tabla de tipos: No olvide que estamos en el fichero `Trans.trg` separado del que contiene las restantes fases de análisis léxico, sintáctico y de ámbito. Además la rutina inicia las referencias `$INT`, `$CHAR` y `$VOID` que serán los árboles/hoja que representen a los tipos básicos. Siguiendo el clásico texto del Dragón de Aho, Hopcroft y Ullman [?] introducimos el tipo `VOID` para asociarlo a aquellos objetos que no tienen tipo (sentencias, etc.).

## 13.7. Análisis de Tipos: Mensajes de Error

### El Problema

Durante esta fase pasamos a asignar un tipo a los nodos del AST. Pueden producirse errores de tipo y es importante que los mensajes de diagnóstico sea precisos. Por ejemplo, dado el programa:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/script> cat -n prueba17.c
 1 int c[20][30], d;
 2
 3 int f(int a, int b) {
 4 return
 5 (a+b)*
 6 d*
 7 c[5];
 8 }

```

queremos producir un mensaje de error adecuado como:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/script> usetypes.pl prueba17.c
Type Error at line 6: Incompatible types with operator '*'

```

indicando así que el error se produce en la multiplicación de la línea 6, ya que `c[5]` es de tipo array.

El problema que tenemos a estas alturas de la construcción de nuestro compilador es que el número de línea asociado con la multiplicación de la línea 6 no figura como atributo del nodo `TIMES`: fue obviado por la directiva `%tree` ya que el terminal `*` fué tratado como un *syntactic token*.

## El método `TERMINAL::save_attributes`

Afortunadamente `Parse::Eyapp` provee un mecanismo para permitir guardar la información residente en un terminal sintáctico en el nodo padre del mismo. Si el programador provee a la clase `TERMINAL` de un método `save_attributes` dicho método será llamado durante la construcción del AST en el momento de la eliminación del terminal:

```
TERMINAL::save_attributes($terminal, $lhs)
```

El primer argumento es la referencia al nodo `TERMINAL` y el segundo al node padre del terminal.

Por tanto, para resolver el problema de conocer el número de línea en el que ocurre un operador, proveeremos a los nodos `TERMINAL` del siguiente método `save_attributes`:

```
620 sub TERMINAL::save_attributes {
621 # $_[0] is a syntactic terminal
622 # $_[1] is the father.
623 push @{$_[1]->{lines}}, $_[0]->[1]; # save the line!
624 }
```

## La Función `insert_method`

La subrutina de mensajes de error de tipo `type_error` la llamaremos proporcionando el mensaje y el número de línea. Sigue un ejemplo de llamada:

```
type_error("Incompatible types with operator '".($_[0]->lexeme)."'", $_[0]->line);
```

El método `line` para este tipo de nodos tiene la forma:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/sub P.*ne/,/^)/p' Types.eypp | cat -n
1 sub PLUS::line {
2 $_[0]->{lines}[0]
3 }
4
5 insert_method(
6 qw{TIMES DIV MINUS ASSIGN GT IF RETURN},
7 'line',
8 \&PLUS::line
9);
```

La función `insert_method` es proveída por `Parse::Eyapp::Base` (versiones de `Parse::Eyapp` posteriores a la 1.099). Recibe como argumentos una lista de clases, el nombre del método (`'line'` en el ejemplo) y la referencia a la función que la implementa. Dotará a cada una de las clases especificadas en la lista (`TIMES`, `DIV` etc.) con métodos `line` implantados a través de la función especificada (`\&PLUS::line`).

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/use/,/^$/p' Types.eypp | cat -n
1 use strict;
2 use Carp;
3 use warnings;
4 use Data::Dumper;
5 use List::MoreUtils qw(firstval);
6 use Simple::Trans;
7 use Parse::Eyapp::Scope qw(:all);
8 use Parse::Eyapp::Base qw(insert_function insert_method);
9 our $VERSION = "0.4";
10
```

## El Lexema Asociado con Una Clase

El método `lexeme` que aparece en la llamada a `type_error` devuelve para un tipo de nodo dado el lexema asociado con ese nodo. Por ejemplo, para los nodos `TIMES` retorna la cadena `*`. Para implantarlo se usan los hashes `lexeme` y `rlexeme`:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '31,64p' Types.eyp | cat -n
1 my %lexeme = (
2 '=' => "ASSIGN",
3 '+' => "PLUS",
4 '-' => "MINUS",
5 '*' => "TIMES",
6 '/' => "DIV",
7 '%' => "MOD",
8 '|' => "OR",
9 '&' => "AND",
10 '{' => "LEFTKEY",
11 '}' => "RIGHTKEY",
12 ',' => "COMMA",
13 ';' => "SEMICOLON",
14 '(' => "LEFTPARENTHESIS",
15 ')' => "RIGHTPARENTHESIS",
16 '[' => "LEFTBRACKET",
17 ']' => "RIGHTBRACKET",
18 '==' => "EQ",
19 '+=' => "PLUSASSIGN",
20 '-=' => "MINUSASSIGN",
21 '*=' => "TIMESASSIGN",
22 '/=' => "DIVASSIGN",
23 '%=' => "MODASSIGN",
24 '!=' => "NE",
25 '<' => "LT",
26 '>' => "GT",
27 '<=' => "LE",
28 '>=' => "GE",
29 '++' => "INC",
30 '--' => "DEC",
31 '**' => "EXP"
32);
33
34 my %rlexeme = reverse %lexeme;
```

Cada una de las clases implicadas es dotada de un método `lexeme` mediante una llamada a `insert_method`:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/sub.*xeme/,/^)/p' Types.eyp | cat -n
1 sub lexeme {
2 return $rlexeme{ref($_[0])} if defined($rlexeme{ref($_[0])});
3 return ref($_[0]);
4 }
5
6 insert_method(
7 # Gives the lexeme for a given operator
8 qw{
9 PLUS TIMES DIV MINUS
```

```

10 GT EQ NE
11 IF
12 IFELSE
13 WHILE
14 VARARRAY VAR ASSIGN
15 FUNCTIONCALL
16 },
17 'lexeme',
18 \&lexeme
19);

```

### Listas y save\_attributes

El método `save_attributes` no se ejecuta automáticamente en las reglas que usan los operadores de listas `+`, `*` y `?`. En tales casos el programador deberá proveer código explícito que salve el número de línea:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> \
 sed -ne '353,363p' Types.eyp | cat -n
 1 Variable:
 2 %name VAR
 3 ID
 4 | %name VARARRAY
 5 $ID ('[' binary ']') <%name INDEXSPEC +>
 6 {
 7 my $self = shift;
 8 my $node = $self->YYBuildAST(@_);
 9 $node->{line} = $ID->[1];
10 return $node;
11 }

```

En el ejemplo se pone explícitamente como atributo `line` del nodo `VARARRAY` el número de línea del terminal `ID`.

## 13.8. Comprobación de Tipos: Las Expresiones

La comprobación de tipos la haremos expandiendo nuestro programa árbol con nuevas transformaciones para cada uno de los tipos de nodos.

### Expresiones Binarias

Las operaciones binarias requieren que sus operandos sean de tipo entero. Si el tipo de alguno de los operandos es `CHAR` haremos una conversión explícita al tipo `INT` (líneas 17-18):

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> sed -ne '66,90p' Trans.trg |
 1 bin: / PLUS
 2 |MINUS
 3 |TIMES
 4 |DIV
 5 |MOD
 6 |GT
 7 |GE
 8 |LE
 9 |EQ
10 |NE
11 |LT

```

```

12 |AND
13 |EXP
14 |OR
15 /($x, $y)
16 => {
17 $x = char2int($_[0], 0);
18 $y = char2int($_[0], 1);
19
20 if (($x->{t} == $INT) and ($y->{t} == $INT)) {
21 $_[0]->{t} = $INT;
22 return 1;
23 }
24 type_error("Incompatible types with operator '" .($_[0]->lexeme)."', $_[0]->line);
25 }

```

## Modificación de la Semántica de las Expresiones Regulares

Observe la expresión regular lineal en las líneas 1-15. La semántica de las expresiones regulares lineales es modificada ligeramente por `Parse::Eyapp::Treeregexp`. Note que no se ha especificado la opción `x`. El compilador de expresiones regulares árbol la insertará por defecto. Tampoco es necesario añadir anclas de frontera de palabra `\b` a los identificadores que aparecen en la expresión regular lineal: de nuevo el compilador de expresiones regulares árbol las insertará.

## Introducción de Nodos de Conversión de Tipos

Las subrutinas de conversión `char2int` e `int2char` se proveen - junto con la subrutina para emitir mensajes de error - como código de apoyo a las expresiones árbol. Recuerde que en un programa `Treeregexp` puede incluir código Perl en cualquier lugar sin mas que aislarlo entre llaves. Dicho código será insertado - respetando el orden de aparición en el fuente - en el módulo generado por `treereg`:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ sed -ne '/^{$/|^}$p' Tr
1 {
2
3 my $types; # reference to the hash containing the type table
4 my ($INT, $CHAR, $VOID);
5
6 sub type_error {
7 my $msg = shift;
8 my $line = shift;
9 die "Type Error at line $line: $msg\n"
10 }
11
12 sub set_types {
13
14 }
15
16 sub char2int {
17 my ($node, $i) = @_;
18
19 my $child = $node->child($i);
20 return $child unless $child->{t} == $CHAR;
21
22 my $coherced = Parse::Eyapp::Node->new('CHAR2INT', sub { $_[0]->{t} = $INT });
23 $coherced->children($child); # Substituting $node(..., $child, ...)
24 $node->child($i, $coherced); # by $node(..., CHAR2INT($child), ...)

```

```

29
30 return $coherced;
31 }
32
33 sub int2char {
34 my ($node, $i) = @_;
35
36 my $child = $node->child($i);
37 return $child unless $child->{t} == $INT;
38
39 my $coherced = Parse::Eyapp::Node->new('INT2CHAR', sub { $_[0]->{t} = $CHAR });
40 $coherced->children($child); # Substituting $node(..., $child, ...)
41 $node->child($i, $coherced); # by $node(..., INT2CHAR($child), ...)
42
43 return $coherced;
44 }
45
46 }

```

Si el hijo  $i$ -ésimo del nodo `$node` es de tipo `$CHAR`, la línea 26 creará un nuevo nodo. Mediante el manejador pasado como segundo argumento a `Parse::Eyapp::Node->new` se dota al nodo del atributo `t`, el cual se establece a `$INT`. El nuevo nodo se interpone entre padre e hijo (líneas 27 y 28). De este modo se facilita la labor de la fase posterior de generación de código, explicitando la necesidad de generar código de conversión.

Si el lenguaje tuviera mas tipos numéricos, `FLOAT`, `DOUBLE` etc. es útil hacer específico en el tipo del nodo que operación binaria se esta realizando cambiando el tipo `PLUS` en `PLUSDOUBLE`, `PLUSINT` etc. según corresponda.

### Expresiones Constantes

La asignación de tipo a los nodos que se corresponden con expresiones constantes es trivial:

```

nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> \
sed -ne '60,61p' Trans.trg | cat -n
 1 inum: INUM($x) => { $_[0]->{t} = $INT }
 2 charconstant: CHARCONSTANT($x) => { $_[0]->{t} = $CHAR }

```

**Ejemplo** Sigue un ejemplo de comprobación de tipos en expresiones:

```

pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba18.c 2
1 int c[20][30], d;
2
3 int f(int a, int b) {
4 if (c[2] > 0)
5 return
6 (a+b)*d*c[1][1];
7 }
Type Error at line 4: Incompatible types with operator '>'

```

## 13.9. Comprobación de Tipos: Indexados

### Guía de Ruta en la Comprobación de Tipos para los Accesos a Arrays

En la comprobación de tipos de una expresión `x[y1][y2]` que referencia a un elemento de un array hay que tener en cuenta:



1. Que el tipo declarado para la variable `x` debe ser array (líneas 4-9 del código que sigue)
2. Que el número de dimensiones utilizadas en la expresión debe ser menor o igual que el declarado para la variable (líneas 11-17).
3. Si `x` es un array de tipo `t` el tipo de `x[2]` será `t`. Por ejemplo, en la llamada `g(x[y1][y2])` la variable `x` deberá haber sido declarada con al menos dos dimensiones. Si `x` fué declarada `int x[10][20][30]` el tipo de la expresión `x[y1][y2]` será `A_30(INT)`.
4. Que los índices usados deben ser de tipo entero. Si fueran de tipo `CHAR` podrán ser ahormados.

### Código de Comprobación de Tipos para los Accesos a Arrays

Sigue el código para la comprobación de tipos de expresiones de referencias a elementos de un array:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> \
sed -ne '97,125p' Trans.trg | cat -n
1 arrays: VARARRAY($x, INDEXSPEC(@y))
2 => {
3
4 my $t = $VARARRAY->{t}; # Type declared for VARARRAY
5 type_error(# Must be an array type
6 " Variable '$x->{attr}[0]' was not declared as array",
7 $VARARRAY->line
8)
9 unless is_array($t);
10
11 my ($declared_dim, $used_dim, $ret_type) = compute_dimensionality($t, @y);
12
13 type_error(
14 " Variable '$x->{attr}[0]' declared with less than $used_dim dimensions",
15 $VARARRAY->line
16)
17 unless $declared_dim >= $used_dim;
18
19 for (0..$#y) { # check that each index is integer. Coerce it if is $CHAR
20 my $ch = char2int($INDEXSPEC, $_);
21
22 type_error("Indices must be integers",$VARARRAY->line)
23 unless ($ch->{t} == $INT);
24 }
25
26 $VARARRAY->{t} = $ret_type;
27
28 return 1;
29 }
```

### Funciones de Ayuda

La función `compute_dimensionality` llamada en la línea 11 se incluye dentro de la sección de funciones de soporte en `Trans.trg`:

```
96 { # support for arrays
97
98 sub compute_dimensionality {
99 my $t = shift;
```

```

100 my $i = 0;
101 my $q;
102 my $used_dim = scalar(@_);
103 for ($q=$t; is_array($q) and ($i < $used_dim); $q=$q->child(0)) {
104 $i++
105 }
106
107 croak "Error checking array type\n" unless defined($q);
108 return ($i, $used_dim, $q);
109 }
110
111 sub is_array {
112 my $type = shift;
113
114 defined($type) && $type =~ /^A_\d+/;
115 }
116
117 sub array_compatible {
118 my ($a1, $a2) = @_;
119
120 return 1 if $a1 == $a2;
121 # int a[10][20] and int b[5][20] are considered compatibles
122 return (is_array($a1) && is_array($a2) && ($a1->child(0) == $a2->child(0)));
123 }
124 }

```

## Ejemplos

El siguiente ejemplo muestra la comprobación de tipos de arrays en funcionamiento:

```
pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba02.c 2
```

```

1 int a,b,e[10];
2
3 g() {}
4
5 int f(char c) {
6 char d;
7 c = 'X';
8 e[d][b] = 'A'+c;
9 {
10 int d;
11 d = a + b;
12 }
13 c = d * 2;
14 return c;
15 }
16

```

Type Error at line 8: Variable 'e' declared with less than 2 dimensions

El siguiente ejemplo muestra la comprobación de tipos de que los índices de arrays deben ser enteros:

```
pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba06.c 2
```

```

1 int a,b,e[10][20];
2
3 g() {}
4

```

```

5 int f(char c) {
6 char d[20];
7 c = 'X';
8 g(e[d], 'A'+c);
9 {
10 int d;
11 d = a + b;
12 }
13 c = d * 2;
14 return c;
15 }
16
Type Error at line 8: Indices must be integers

```

### 13.10. Comprobación de Tipos: Sentencias de Control

Las sentencias de control como IF y WHILE requieren un tipo entero en la condición de control:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/^control/,189p' Trans.trg | cat -n
1 control: /IF|IFELSE|WHILE/:con($bool)
2 => {
3 $bool = char2int($con, 0);
4
5 type_error("Condition must have integer type!", $bool->line)
6 unless $bool->{t} == $INT;
7
8 $con->{t} = $VOID;
9
10 return 1;
11 }

```

Veamos un ejemplo de uso:

```

l@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba19.c 2
1 int c[20][30], d;
2
3 int f(int a, int b) {
4 if (c[2])
5 return
6 (a+b)*d*c[1][1];
7 }
Type Error at line 4: Condition must have integer type!
pl@nereida:~/Lbook/code/Simple-Types/script$

```

### 13.11. Comprobación de Tipos: Sentencias de Asignación

La comprobación de tipos en las asignaciones se divide en tres fases:

1. Coherción: líneas 10 y 11
2. Comprobación de la igualdad de los tipos (DAGs) de ambos lados de la igualdad (líneas 13 y 14)

3. Comprobación de que el tipo del lado izquierdo de la asignación está entre los permitidos por el lenguaje. En ANSI C se permite también la asignación de estructuras pero estas no están presentes en Simple C.
4. Si el lenguaje permitiera expresiones mas complejas en el lado izquierdo serían necesarias comprobaciones de que la expresión en la izquierda es 'asignable' (es un lvalue). Por ejemplo en Kernighan y Ritchie C (KR C) [?] la asignación `f(2,x) = 4` no es correcta pero si podría serlo `*f(2,x) = 4`. En KR C las reglas de producción que definen la parte izquierda de una asignación (lvalue) son mas complejas que las de SimpleC, especialmente por la presencia de punteros:

| Resumen de la Gramática de Kernighan y Ritchie C (KR C) en Parse::Eyapp                                                                                                                                                                                                                                                                                                                                            |                                                                                                                                                                                                                                                                                                                                                                      |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Expresiones                                                                                                                                                                                                                                                                                                                                                                                                        | Declaraciones                                                                                                                                                                                                                                                                                                                                                        |
| <pre> lvalue:     identifier     primary '[' exp ']'     lvalue '.' identifier     primary '-&gt;' identifier     '*' exp     '(' lvalue ')' primary:     identifier     constant     string     (exp)     primary '(' exp '&lt;*, ','&gt; ')'     primary '[' exp ']'     lvalue '.' identifier     primary '-&gt;' identifier exp:     primary     '*' exp     '&amp;' exp     lvalue ASGNOP exp     .....</pre> | <pre> fund:     t? dec '(' pr '&lt;*, ','&gt; ')' fb declaration:     scty * (dec ini) '&lt;*, '&gt;' scty:     scope     t scope:     tdef   ... t:     char     int     double     tdef_name     ..... dec:     identifier     '(' dec ')'     '*' dec     dec '(' ')'     dec '[' cexp '?' ']' ini:     '=' exp     '=' '{' exp '&lt;*, ','&gt; ',' '?' '}'</pre> |

5. En la línea 20 modificamos el tipo del nodo: de ser ASSIGN pasa a ser ASSIGNINT o ASSIGNCHAR según sea el tipo de la asignación. De este modo eliminamos la sobrecarga semántica del nodo y hacemos explícito en el nodo el tipo de asignación. Poco a poco nos vamos acercando a niveles mas bajos de programación.

Sigue el código de comprobación para las asignaciones:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/^assign/,180p' Trans.trg | cat -n
1 assign: /ASSIGN
2 |PLUSASSIGN
3 |MINUSASSIGN
4 |TIMESASSIGN
5 |DIVASSIGN
6 |MODASSIGN
```

```

7 /:asgn($lvalue, $exp)
8 => {
9 my $lt = $lvalue->{t};
10 $exp = char2int($asgn, 1) if $lt == $INT;
11 $exp = int2char($asgn, 1) if $lt == $CHAR;
12
13 type_error("Incompatible types in assignment!", $asgn->line)
14 unless ($lt == $exp->{t});
15
16 type_error("The C language does not allow assignments to non-scalar types!", $asgn->
17 unless ($lt == $INT) or ($lt == $CHAR); # Structs will also be allowed
18
19 # Assignments are expressions. Its type is the type of the lhs or the rhs
20 $asgn->{t} = $lt;
21
22 # Make explicit the type of assignment, i.e. s/PLUSASSIGN/PLUSASSIGNINT/
23 $asgn->type(ref($asgn).ref($lt));
24
25 return 1;
26 }

```

Veamos un ejemplo de comprobación de tipos en asignaciones:

```
pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba12.c 2
```

```
1 int a,b[10][20],e[10][20];
```

```
2
```

```
3 int f(char c) {
```

```
4 e[5] = b[5];
```

```
5 a = b[1][2];
```

```
6 }
```

```
7
```

```
Type Error at line 4: The C language does not allow assignments to non-scalar types!
```

## 13.12. Comprobación de Tipos: Llamadas a Funciones

### Guía de Ruta para la Comprobación de Tipos de las LLamadas

La comprobación de tipos en una llamada a función conlleva los siguientes pasos:

1. Líneas 4 y 5: comprobamos que la variable que se usa en la llamada fue efectivamente declarada como función
2. Línea 14: El número de argumentos y el número de parámetros declarados debe coincidir
3. En las líneas de la 19 a la 35 analizamos la compatibilidad de tipo de cada uno de los argumentos con su correspondiente parámetro:
  - a) En la línea 21 se obtiene el nodo correspondiente al argumento. Si es necesario se harán las coerciones (líneas 22 y 23)
  - b) En la línea 25 comprobamos la compatibilidad de tipos. En el caso de los arrays la compatibilidad no se limita a la igualdad de referencias: consideraremos compatibles arrays que difieren en la primera dimensión. Así las declaraciones `int a[10][20]` y `int b[5][20]` se considerarán compatibles pero `int c[10][10]` y `int d[10][20]` no.

## La forma del Árbol de las LLamadas

La siguiente sesión con el depurador ilustra la forma del árbol de análisis en las llamadas:

```
pl@nereida:~/Lbook/code/Simple-Types/script$ perl -wd usetypes.pl prueba23.c 2
main::(usetypes.pl:5): my $filename = shift || die "Usage:\n$0 file.c\n";
```

```
DB<1> f Types.eyp
1 2 3 4 5 6 7 #line 8 "Types.eyp"
8
9: use strict;
10: use Carp;
DB<2> c 598
1 int f(char a[10], int b) {
2 return a[5];
3 }
4
5 int h(int x) {
6 return x*2;
7 }
8
9 int g() {
10 char x[5];
11 int y[19][30];
12 f(x,h(y[1][1]));
13 }
Simple::Types::compile(Types.eyp:598):
598: set_types($t); # Init basic types
```

Ejecutamos hasta la línea 598, justo después de terminar el análisis de ámbito y antes de comenzar con la fase de comprobación de tipos.

A continuación mostramos una versión compacta del árbol. Para ello suprimimos la salida de notas a pie de página "vaciando" el método `footnotes`:

```
DB<3> insert_method(qw{PROGRAM FUNCTION}, 'footnote', sub {}) # Avoid footnotes
DB<4> p $t->str
PROGRAM(
 FUNCTION[f] (RETURN (VARARRAY (TERMINAL[a:2], INDEXSPEC (INUM (TERMINAL[5:2]))))),
 FUNCTION[h] (RETURN (TIMES (VAR (TERMINAL[x:6]), INUM (TERMINAL[2:6])))),
 FUNCTION[g] (
 FUNCTIONCALL(
 TERMINAL[f:12],
 ARGLIST(
 VAR (TERMINAL[x:12]),
 FUNCTIONCALL(
 TERMINAL[h:12],
 ARGLIST (VARARRAY (TERMINAL[y:12], INDEXSPEC (INUM (TERMINAL[1:12])), INUM (TERMINAL[1:12]))))
)
)
)
)
)
DB<6> p $t->descendant(".2.0")->str
FUNCTIONCALL(
 TERMINAL[f:12],
 ARGLIST(
 VAR (TERMINAL[x:12]),
 FUNCTIONCALL(
 TERMINAL[h:12],
 ARGLIST (VARARRAY (TERMINAL[y:12], INDEXSPEC (INUM (TERMINAL[1:12])), INUM (TERMINAL[1:12]))))
)
)
)
DB<7> p $t->descendant(".2.0")->{t}->str
```

```
F(X_2(A_10(CHAR),INT),INT)
```

La sesión ilustra como la función `insert_method` puede ser utilizada para introducir métodos en las clases influyendo en la presentación del árbol. Veamos un segundo ejemplo en el que eliminamos la presentación de la información asociada con los nodos:

```
DB<8> insert_method(qw{FUNCTION TERMINAL}, 'info', sub {})
DB<9> x $t->str
0 'PROGRAM(
 FUNCTION(RETURN(VARARRAY(TERMINAL,INDEXSPEC(INUM(TERMINAL))))),
 FUNCTION(RETURN(TIMES(VAR(TERMINAL),INUM(TERMINAL)))),
 FUNCTION(
 FUNCTIONCALL(TERMINAL,ARGLIST(VAR(TERMINAL),i
 FUNCTIONCALL(
 TERMINAL,
 ARGLIST(VARARRAY(TERMINAL,INDEXSPEC(INUM(TERMINAL),INUM(TERMINAL)))))))))',
```

### Código de Comprobación de las LLamadas

Sigue el código de comprobación de tipos de las llamadas:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '/^func/,234p' Trans.trg | cat -n
1 functioncall: FUNCTIONCALL($f, ARGLIST)
2 => {
3 # Before type checking attribute "t" has the declaration of $f
4 my $ftype = $FUNCTIONCALL->{t};
5
6 type_error(" Variable '". $f->value.'" was not declared as function", $f->line)
7 unless $ftype->isa("F");
8
9 my @partypes = $ftype->child(0)->children;
10
11 my @args = $ARGLIST->children; # actual arguments
12 my $numargs = @args; # Number of actual arguments
13 my $numpar = @partypes; # Number of declared parameters
14
15 # Check number of args
16 type_error("Function '". $f->value.'" called with $numargs args expected $numpar",$f->
17 if ($numargs != $numpar);
18
19 # Check type compatibility between args
20 # Do type coercion if needed
21 for (0..$#args) {
22 my $pt = shift @partypes;
23 my $ch = $ARGLIST->child($_);
24 $ch = char2int($ARGLIST, $_) if $pt == $INT;
25 $ch = int2char($ARGLIST, $_) if $pt == $CHAR;
26
27 my $cht = $ch->{t};
28 unless (array_compatible($cht, $pt)) {
29 type_error(
30 "Type of argument " .($_+1)." in call to " . $f->value." differs from expected", $f->
31 $f->line
```

```

32)
33 }
34 }
35
36 # Now attribute "t" has the type of the node
37 $FUNCTIONCALL->{t} = $ftype->child(1);
38 return 1;
39 }

```

### Arrays que Difieren en la Primera Dimensión

La función `array_compatible` sigue la costumbre C de considerar compatibles arrays que difieren en la primera dimensión. Véase la conducta de `gcc` al respecto:

```

pl@nereida:~/Lbook/code/Simple-Types/script$ cat -n ArrayChecking.c
 1 #include <stdio.h>
 2
 3 void f(int a[3][3]) {
 4 a[1][2] = 5;
 5 }
 6
 7 main() {
 8 int b[2][3];
 9 f(b);
10 printf("%d\n",b[1][2]);
11 }
pl@nereida:~/Lbook/code/Simple-Types/script$ gcc ArrayChecking.c
pl@nereida:~/Lbook/code/Simple-Types/script$ a.out
5

```

Recordemos el código de la función `array_compatible`:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ sed -ne '96,124p' Trans.t
 1 { # support for arrays
 2
 3 sub compute_dimensionality {
..
14 }
15
16 sub is_array {
..
20 }
21
22 sub array_compatible {
23 my ($a1, $a2) = @_;
24
25 return 1 if $a1 == $a2;
26 # int a[10][20] and int b[5][20] are considered compatibles
27 return (is_array($a1) && is_array($a2) && ($a1->child(0) == $a2->child(0)));
28 }
29 }

```

### Ejemplos

```

pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba07.c 2

```



```

1 int a,b,e[10][20];
2
3 g() {}
4
5 int f(char c) {
6 char d;
7 c = 'X';
8 g(e[d], 'A'+c);
9 if (c > 0) {
10 int d;
11 d = a + b;
12 }
13 c = d * 2;
14 return c;
15 }
16

```

Type Error at line 8: Function 'g' called with 2 args expected 0

```

pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba22.c 2

```

```

1 int f(char a[10], int b) {
2 return a[5];
3 }
4
5 int h(int x) {
6 return x*2;
7 }
8
9 int g() {
10 char x[5];
11 int y[19][30];
12 f(x,h(y));
13 }

```

Type Error at line 12: Type of argument 1 in call to h differs from expected

### 13.13. Comprobación de Tipos: Sentencia RETURN

#### La Forma Habitual

Veamos la forma del árbol en una sentencia de retorno:

```

pl@nereida:~/Lbook/code/Simple-Types/script$ perl -wd usetypes.pl prueba16.c 2
main::(usetypes.pl:5): my $filename = shift || die "Usage:\n$0 file.c\n";
DB<1> f Types.eyp
1 2 3 4 5 6 7 #line 8 "Types.eyp"
8
9: use strict;
10: use Carp;
DB<2> c 598
1 int f() {
2 int a[30];
3
4 return a;
5 }
Simple::Types::compile(Types.eyp:598):
598: set_types($t); # Init basic types

```

```
DB<3> insert_method(qw{PROGRAM FUNCTION}, 'footnote', sub {})
DB<4> p $t->str
PROGRAM(FUNCTION[f] (RETURN(VAR(TERMINAL[a:4]))))
```

## Código de Comprobación de Tipos para Sentencias de Retorno

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '267,282p' Trans.trg | cat -n
 1 returntype: RETURN($ch)
 2 => {
 3 my $rt = $RETURN->{t};
 4
 5 $ch = char2int($RETURN, 0) if $rt == $INT;
 6 $ch = int2char($RETURN, 0) if $rt == $CHAR;
 7
 8 type_error("Type error in return statement", $ch->line)
 9 unless ($rt == $ch->{t});
10
11 # $RETURN->{t} has already the correct type
12
13 $RETURN->type(ref($RETURN).ref($rt));
14
15 return 1;
```

El análisis de ámbito de las sentencias de retorno se hizo de forma ligeramente distinta al resto. El siguiente fragmento de código figura en la subrutina `compile` justo después del resto del análisis de ámbito:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple$ \
 sed -ne '585,595p' Types.eyp | cat -n
 1 # Scope Analysis: Return-Function
 2 our $retscope; # retscope: /FUNCTION|RETURN/
 3 my @returns = $retscope->m($t);
 4 for (@returns) {
 5 my $node = $_->node;
 6 if (ref($node) eq 'RETURN') {
 7 my $function = $_->father->node;
 8 $node->{function} = $function;
 9 $node->{t} = $function->{t}->child(1);
10 }
11 }
```

## Ejemplos

```
pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba16.c 2
1 int f() {
2 int a[30];
3
4 return a;
5 }
Type Error at line 4: Type error in return statement
```

El retorno vacío por parte de una función no se considera un error en C:

```
pl@nereida:~/Lbook/code/Simple-Types/script$ usetypes.pl prueba26.c 2 | head -12
1 int f() {
2 int a[30];
3
4 return;
5 }
```

```
PROGRAM^{0}{
 FUNCTION[f]^{1}{
 EMPTYRETURN
 }
} # PROGRAM

```

Véase la conducta de gcc ante el mismo programa:

```
pl@nereida:~/Lbook/code/Simple-Types/script$ gcc -c prueba26.c
pl@nereida:~/Lbook/code/Simple-Types/script$ ls -ltr | tail -1
-rw-r--r-- 1 pl users 690 2008-01-10 13:32 prueba26.o
```

## Usando m

En esta sección seguiremos una aproximación diferente. En este caso queremos maximizar nuestro conocimiento

Para comprobar la corrección del tipo del valor retornado vamos a usar el método `m`. La razón para hacerlo es familiarizarle con el uso de `m`.

Usaremos una transformación árbol `bind_ret2function` a la que llamaremos (línea 33) desde `compile` después de haber verificado el resto de la comprobación de tipos (línea 30). Podemos post-poner este análisis ya que el padre de un nodo `RETURN` es un nodo `FUNCTION` y por tanto no hay nodos cuya comprobación de tipos dependa de la computación de tipo de un nodo `RETURN`.

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> \
 sed -ne '519,554p' Types.eypp | cat -n
1 sub compile {
2 my($self)=shift;
3
4
29
30 $t->bud(@typecheck);
31
32 our $bind_ret2function;
33 my @FUNCTIONS = $bind_ret2function->m($t);
34
35 return $t;
36 }
```

`bind_ret2function` realiza la comprobación de tipos en los siguientes pasos:

1. La expresión regular árbol casa con cualquier nodo de definición de función. Para cada uno de esos nodos se procede como sigue:
2. Se buscan (línea 4) todas las sentencias `RETURN` en la función. Recuerde que `Parse::Eyapp::Node::m` retorna nodos `Match` y que los nodos del AST se obtienen usando el método `node` sobre el nodo `Match` (línea 5)
3. Ponemos como atributo de la función las referencias a los nodos `RETURN` encontrados (línea 8)

4. Para cada nodo de tipo **RETURN** encontrado:

- a) Ponemos como atributo **function** del nodo la función que le anida (línea 15)
- b) Se realizan las coerciones que fueran necesarias (líneas 17-19)
- c) Se emite un mensaje de error si el tipo de la expresión evaluada ( $\$exp \rightarrow \{t\}$ ) no coincide con el declarado en la función ( $\$return\_type$ ) (líneas 21-22)
- d) En la línea 25 cambiamos la clase del nodo de **RETURN** a **RETURNINT** o **RETURNCHAR** según sea el tipo retornado. Al igual que en las asignaciones eliminamos la sobrecarga semántica del nodo acercándonos a niveles mas bajos de programación.

Sigue el código. La línea 1 muestra una simple expresión regular árbol que casa con aquellas sentencias **RETURN** en las que se ha explicitado una expresión de retorno y que es usada por **bind\_ret2function**.

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> \
sed -ne '196,225p' Trans.trg | cat -n
1 /* TIMTOWTDI when MOPping */
2 return: RETURN(.)
3 bind_ret2function: FUNCTION
4 => {
5 my @RETURNS = $return->m($FUNCTION);
6 @RETURNS = map { $_->node } @RETURNS;
7
8 # Set "returns" attribute for the FUNCTION node
9 $FUNCTION->{returns} = \@RETURNS;
10
11 my $exp;
12 my $return_type = $FUNCTION->{t}->child(1);
13 for (@RETURNS) {
14
15 # Set "function" attribute for each RETURN node
16 $_->{function} = $FUNCTION;
17
18 #always char-int conversion
19 $exp = char2int($_, 0) if $return_type == $INT;
20 $exp = int2char($_, 0) if $return_type == $CHAR;
21
22 type_error("Returned type does not match function declaration",
23 $_->line)
24 unless $exp->{t} == $return_type;
25 $_->type("RETURN".ref($return_type));
26
27 }
28
29 return 1;
30 }
```

Obsérvese que la solución no optimiza la homogeneidad ni la generalidad ni el mantenimiento. Por ello, en condiciones normales de construcción de un compilador - como es el caso de la práctica 13.16 - es aconsejable abordar el análisis de tipos de la sentencia **RETURN** con la misma aproximación seguida para el resto de los nodos.

### 13.14. Comprobación de Tipos: Sentencias sin tipo

Usamos el tipo **\$VOID** como tipo de aquellos constructos del lenguaje que carecen de tipo:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/lib/Simple> \
sed -ne '67p' Trans.trg
statements: /STATEMENTS|PROGRAM|BREAK|CONTINUE/ => { $_[0]->{t} = $VOID }
```

### 13.15. Ejemplo de Árbol Decorado

Consideremos el programa de ejemplo:

```
nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-Types/script> cat -n prueba23.c
1 int f(char a[10], int b) {
2 return a[5];
3 }
4
5 int h(int x) {
6 return x*2;
7 }
8
9 int g() {
10 char x[5];
11 int y[19][30];
12 f(x,h(y[1][1]));
13 }
```

La siguiente tabla muestra (parcialmente) el árbol decorado resultante después de la fase de análisis de tipos:

| Funciones f y h                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                | Función g                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> PROGRAM^{0}{(   FUNCTION[f](     RETURNINT( # f(char a[10],       CHAR2INT( # return a[5];       VARARRAY( # a[5]         TERMINAL[a:2],         INDEXSPEC(           INUM(             TERMINAL[5:2]           )         ) # INDEXSPEC       ) # VARARRAY     ) # CHAR2INT   ) # RETURNINT ) # FUNCTION, FUNCTION[h]( # int h(int x)   RETURNINT( # return x*2     TIMES[INT:6](       VAR(         TERMINAL[x:6]       ),       INUM(         TERMINAL[2:6]       )     ) # TIMES   ) # RETURNINT ) # FUNCTION, </pre> | <pre> FUNCTION[g](   FUNCTIONCALL( # f(x,h(y[1][1]))     TERMINAL[f:12],     ARGLIST(       VAR( # char x[5]         TERMINAL[x:12]       ),       FUNCTIONCALL( # h(y[1][1])         TERMINAL[h:12],         ARGLIST(           VARARRAY( # y[1][1]             TERMINAL[y:12],             INDEXSPEC(               INUM(                 TERMINAL[1:12]               ),               INUM(                 TERMINAL[1:12]               )             ) # INDEXSPEC           ) # VARARRAY         ) # ARGLIST       ) # FUNCTIONCALL     ) # ARGLIST   ) # FUNCTIONCALL ) # FUNCTION ) # PROGRAM </pre> |

### 13.16. Práctica: Análisis de Tipos en Simple C

Realice el analizador de tipos para Simple C de acuerdo a lo explicado en las secciones anteriores. Extienda el análisis de tipos a la extensión del lenguaje con `structs`. Si lo desea puede usar una aproximación OOP estableciendo un método `checktypes` para cada clase/nodo. Para realizar esta alternativa es conveniente definir la jerarquía de clases adecuada.

### 13.17. Práctica: Análisis de Tipos en Simple C

Extienda el análisis de tipos de la práctica anterior para incluir punteros. (repase el resumen de la gramática de KR C en la página 725).

### 13.18. Práctica: Análisis de Tipos en Simple C con Gramáticas Atribuidas

Reescriba el analizador de tipos para Simple C utilizando `Language::AttributeGrammars`. Repase la sección 9.5 (*Usando Language::AttributeGrammars con Parse::Eyapp*).

Es conveniente que las acciones sean aisladas en subrutinas para evitar que `Language::AttributeGrammars` se confunda. Sigue un ejemplo de como podría hacerse:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-AttributeGrammar/lib$ \
```

```

 sed -ne '585,631p' Trans_Scheme.eyp
my $attgram = new Language::AttributeGrammar <<'EOG';

Tipos Basicos
INUM: $/.t = { Trans_Scheme::inum }
CHARCONSTANT: $/.t = { Trans_Scheme::char_constant }
Variables escalares y arrays
VAR: $/.t = { Trans_Scheme::var($/) }
Expresiones binarias
PLUS: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
MINUS: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
TIMES: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
DIV: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
MOD: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
GT: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
GE: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
LE: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
EQ: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
NE: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
LT: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
AND: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
OR: $/.t = { Trans_Scheme::bin($<0>.t, $<1>.t, $/) }
Sentencias de control
IF: $/.t = { Trans_Scheme::condition($<0>.t, $/, $<1>.t) }
IFELSE: $/.t = { Trans_Scheme::condition($<0>.t, $/, $<1>.t, $<2>.t) }
WHILE: $/.t = { Trans_Scheme::condition($<0>.t, $/, $<1>.t) }
Asignaciones
ASSIGN: $/.t = { Trans_Scheme::assign($<0>.t, $<1>.t, $/) }
PLUSASSIGN: $/.t = { Trans_Scheme::assign($<0>.t, $<1>.t, $/) }
MINUSASSIGN: $/.t = { Trans_Scheme::assign($<0>.t, $<1>.t, $/) }
TIMESASSIGN: $/.t = { Trans_Scheme::assign($<0>.t, $<1>.t, $/) }
DIVASSIGN: $/.t = { Trans_Scheme::assign($<0>.t, $<1>.t, $/) }
MODASSIGN: $/.t = { Trans_Scheme::assign($<0>.t, $<1>.t, $/) }
Llamadas a funciones
FUNCTIONCALL: $/.t = { Trans_Scheme::functioncall($/) }
Return
RETURN: $/.t = { Trans_Scheme::return($<0>.t, $/) }
Otros nodos
LABEL: $/.t = { Trans_Scheme::void }
GOTO: $/.t = { Trans_Scheme::void }
STATEMENTS: $/.t = { Trans_Scheme::void }
BLOCK: $/.t = { Trans_Scheme::void }
BREAK: $/.t = { Trans_Scheme::void }
CONTINUE: $/.t = { Trans_Scheme::void }
FUNCTION: $/.t = { Trans_Scheme::void }
PROGRAM: $/.t = { Trans_Scheme::void }

EOG

```

Las funciones de soporte son similares a las que hemos usado en Trans.trg:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Simple-AttributeGrammar/lib$ sed -ne '716,721p' Tr
sub return {
 my ($t, $fathernode) = @_;

```

```

my $child = $fathernode->child(0);
return $t if ($types->{$fathernode->{returntype}} == $t);
type_error("Type error in return statement", $fathernode->line);
}

```

### 13.19. Práctica: Sobrecarga de Funciones en Simple C

Introduzca sobrecarga de funciones en Simple C y modifique el análisis de ámbito y de tipos. Repase las secciones 13.1 y 13.2.

### 13.20. Análisis de Tipos de Funciones Polimorfas

El uso de variables en las expresiones de tipo nos permite hablar de tipos desconocidos. Así un tipo como:

$$F(\text{LIST}(\text{TYPEVAR}::\text{ALPHA}), \text{TYPEVAR}::\text{ALPHA})$$

nos permite hablar de una función que recibe listas de objetos de un tipo desconocido  $\text{TYPEVAR}::\text{ALPHA}$  y retorna objetos de ese mismo tipo  $\text{TYPEVAR}::\text{ALPHA}$ . Una variable de tipo representa un objeto no declarado o no completamente declarado. La *Inferencia de Tipos* es el problema de determinar el tipo de los nodos de uso en el árbol sintáctico abstracto cuando no se conocen total o parcialmente los tipos de los objetos en el programa fuente.

#### 13.20.1. Un Lenguaje con Funciones Polimorfas

En esta sección introducimos un pequeño lenguaje que permite el polimorfismo paramétrico.

##### El Cuerpo

Los programas generados por esta gramática consisten en secuencias de declaraciones seguidas de secuencias de las expresiones cuyos tipos vamos a inferir. Por ejemplo:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ \
cat -n prueba01.ply
1 first : list(ALPHA) -> ALPHA;
2 q : list(list(int))
3 {
4 first(first(q))
5 }

```

Las declaraciones admiten variables de tipo. Las variables de tipo son identificadores formados por letras en mayúsculas. Los identificadores estan formados por minúsculas. La gramática que define el lenguaje es:

```

p: d <+ ';'> '{' e <%name EXPS + ';'> '}'
;

d: id ':' t
;

t: INT
 | STRING
 | TYPEVAR
 | LIST '(' t ')'
 | t '*' t
 | t '->' t

```



```

 | '(' t ')'
;

e: id '(' optargs ')'
 | id
;

optargs:
 /* empty */
 | args
;

args: e
 | args ',' e
;

id: ID
;

```

Para las expresiones nos limitamos a construir el árbol usando bypass automático. Las expresiones se reducen al uso de identificadores y llamadas a función:

```

47 p:
48 d <+ ';' '>' '{' e <%name EXPS + ';' '>.expressions '}'
49 {
50 $expressions->{symboltable} = \"%st;
51 $expressions
52 }
53 ;
54
55 d:
56 $id ':' $t
57 {
58 $st{$id->key} = { ts => $t };
59 }
60 ;
61
62 t:
63 INT
64 {
65 'INT';
66 }
67 | STRING
68 {
69 'STRING';
70 }
71 | TYPEVAR
72 {
73 "Parse::Eyapp::Node::TYPEVAR::$_[1]->[0]"
74 }
75 | LIST '(' $t ')'
76 {
77 "LIST($t)"
78 }

```

```

79 | t.left '*' t.right
80 {
81 "X($left,$right)"
82 }
83 | t.domain '->' t.image
84 {
85 "F($domain,$image)"
86 }
87 | '(' $t ')'
88 {
89 $t;
90 }
91 ;
92
93 e:
94 %name FUNCTIONCALL
95 id '(' optargs ')'
96 | id
97 ;
98
99 optargs:
100 /* empty */
101 | args
102 ;
103
104 args:
105 e
106 | %name ARGS
107 args ',' e
108 ;
109
110 id:
111 %name ID
112 ID
113 ;

```

Para cada declaración se construye un término que describe el tipo y se almacena en la tabla de símbolos %st. Por ejemplo, las expresiones de tipo construidas a partir de las declaraciones del programa anterior son:

```

first type: F(LIST(Parse::Eyapp::Node::TYPEVAR::ALPHA),Parse::Eyapp::Node::TYPEVAR::ALPHA)
q type: LIST(LIST(INT))

```

Usaremos el espacio de nombres `Parse::Eyapp::Node::TYPEVAR` para representar los nodos variable. El árbol resultante para el programa anterior es:

|                                                                                                  |                                                                                                                                   |
|--------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------|
| <pre> 1 first : list(ALPHA) -&gt; ALPHA; 2 q : list(list(int)) 3 { 4   first(first(q)) 5 }</pre> | <pre> EXPS(   FUNCTIONCALL(     ID[first],     FUNCTIONCALL(       ID[first],       ID[q]     )   ) # FUNCTIONCALL ) # EXPS</pre> |
|--------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------|

Obsérvese que no hay asignación ni operaciones binarias en el lenguaje (aunque estas últimas pueden ser emuladas mediante funciones `add`, `times`, etc.). Por ello la comprobación de tipos se centra en las llamadas a función.

## La Cola

En la cola figuran la subrutina de análisis léxico y la de tratamiento de errores. Son una adaptación de las presentadas en las secciones anteriores para `Simple::Types`.

El análisis de tipo para este lenguaje con variables de tipo lo haremos usando expresiones regulares árbol (archivo `Trans.trg`) así como un módulo que implanta el algoritmo de unificación y que hemos llamado `Aho::Unify`. En concreto tenemos los siguientes ficheros en la librería:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho$ ls -l
total 112
-rw-r--r-- 1 pl users 639 2008-01-07 13:01 Makefile
-rw-r--r-- 1 pl users 3837 2008-01-11 16:47 Polymorphism.eyp
-rw-r--r-- 1 pl users 945 2008-01-11 12:20 Trans.trg
-rw-r--r-- 1 pl users 3212 2008-01-11 12:57 Unify.pm
```

Los ficheros son compilados con los comandos:

```

treereg -nonumbers -m Aho::Polymorphism Trans.trg
eyapp -m Aho::Polymorphism Polymorphism.eyp
```

La llamada `$t = $self->YYParse` de la línea 192 realiza el análisis sintáctico y el análisis de ámbito. El análisis de ámbito es trivial ya que hay un sólo ámbito.

```

115 %%
...
186 sub compile {
187 my($self)=shift;
188
189 my ($t);
190
191 $self->YYData->{INPUT} = $_[0];
192
193 $t = $self->YYParse(yylex => \&_Lexer, yyerror => \&_Error,
194 #yydebug => 0x1F
195);
196
197 my @typecheck = (
198 our $functioncall,
199 our $tuple,
200 our $id,
201);
```

```

202
203 Aho::Unify->set(
204 key => 'set',
205 isvar => sub { $_[0] =~ /^Parse::Eyapp::Node::TYPEVAR/ },
206 samebasic => sub {
207 my ($s, $t) = @_;
208
209 return (((ref($s) eq 'INT') || (ref($s) eq 'STRING')) && (ref($s) eq ref($t)));
210 },
211 debug => $debug,
212);
213
214 $t->bud(@typecheck);
215
216 return $t;
217 }
218
219 sub ID::key {
220 return $_[0]->{attr}[0]
221 }
222
223 insert_method('ID', 'info', \&ID::key);

```

Las líneas de la 197 a la 212 se encargan de la inferencia de tipos. Volveremos sobre ellas mas adelante. Por ahora observe que hay tres transformaciones:

- `$functioncall` se encarga de la inferencia y comprobación para las llamadas,
- `$tuple` del cálculo de tipos para las expresiones `args`, `exp` que aparecen en los argumentos. Estas expresiones dan lugar a árboles de la forma `ARGS($x, $y)`. Por ejemplo, en el siguiente programa:

```

f : ALPHA * BETA * GAMMA -> ALPHA;
q : int;
r : int;
s : list(int)
{
 f(q,r,s)
}

```

El árbol construido para la expresión es:

```
EXPS(FUNCTIONCALL(ID[f],ARGS(ARGS(ID[q],ID[r]),ID[s]))))
```

mientras que los tipos son:

```

r type: INT
q type: INT
s type: LIST(INT)
f type: F(
 X(
 X(
 Parse::Eyapp::Node::TYPEVAR::ALPHA,
 Parse::Eyapp::Node::TYPEVAR::BETA

```

```

),
 Parse::Eyapp::Node::TYPEVAR::GAMMA
),
 Parse::Eyapp::Node::TYPEVAR::ALPHA
)

```

- La tercera transformación en `Trans.trg` es `$id`, la cual se encarga del cálculo de los tipos para los nodos ID. Como las anteriores, esta transformación decorará el nodo con un atributo `t` denotando el DAG que describe al tipo. En este caso se tomará el tipo de la tabla de símbolos y se crearán nuevas instancias de las variables de tipo si las hubiera.

El comportamiento del algoritmo de unificación depende de un conjunto de parámetros que es inicializado mediante la llamada a la función `set` del módulo `Aho::Unify`.

La llamada `$t->bud(@typecheck)` aplica las transformaciones en un recorrido abajo-arriba infiriendo los tipos y decorando los nodos.

## La Cabecera

La cabecera muestra la carga de los módulos para la unificación (`Aho::Unify`) y para la inferencia y comprobación de tipos (`Aho::Trans`).

## La directiva `%strict`

Por defecto `Parse::Eyapp` decide que aquellos símbolos que no figuran en la parte izquierda de una regla son terminales (tokens). La directiva `%strict` utilizada en la línea 7 fuerza la declaración de todos los terminales. Cuando se usa la directiva `%strict` el compilador `eyapp` emite un mensaje de advertencia si detecta algún símbolo que no figura en la parte izquierda de alguna regla y no ha sido explícitamente declarado como terminal (mediante directivas como `%token`, `%left`, etc.).

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho$ \
 cat -n Polymorphism.eyp

```

```

1 /*
2 File: Aho/Polymorphism.eyp
3 To build it, Do make or:
4 eyapp -m Aho::Polymorphism Polymorphism.eyp;
5 treereg -m Aho::Polymorphism Trans.trg
6 */
7 %strict
8 %{
9 use strict;
10 use Carp;
11 use warnings;
12 use Aho::Unify qw(:all);
13 use Aho::Trans;
14 use Data::Dumper;
15 our $VERSION = "0.3";
16
17 my $debug = 1;
18 my %reserved = (
19 int => "INT",
20 string => "STRING",
21 list => "LIST",
22);
23
24 my %lexeme = (
25 ',' => "COMMA",

```

```

26 ';' => "SEMICOLON",
27 ':' => "COLON",
28 '*' => "TIMES",
29 '(' => "LEFTPARENTHESIS",
30 ')' => "RIGHTPARENTHESIS",
31 '{' => "LEFTCURLY",
32 '}' => "RIGHTCURLY",
33 '->' => "FUNCTION",
34);
35 my ($tokenbegin, $tokenend) = (1, 1);
36 our %st; # Symbol table
37 %}
38
39 %tree bypass
40
41 %token ID TYPEVAR STRING INT LIST
42
43 %right '->'
44 %left '*'
45
46 %%

```

### Ambigüedades

La regla de producción  $t \rightarrow t \text{ '*' } t$  es obviamente ambigua. La directiva `%left '*'` permite deshacer la ambigüedad. Lo mismo ocurre con la regla para las funciones  $t \rightarrow t \text{ '->' } t$ . En este caso nos decidimos por una asociatividad a derechas, de modo que una declaración como `int -> int -> int` es interpretado como una función que recibe enteros y devuelve funciones. Una tercera fuente de ambigüedad se produce en expresiones como:

$$\text{int} * \text{int} \rightarrow \text{int}$$

que puede ser interpretado como  $\text{int} * (\text{int} \rightarrow \text{int})$  o bien  $(\text{int} * \text{int}) \rightarrow \text{int}$ . Al darle mayor prioridad al `*` nos decidimos por la segunda interpretación.

### 13.20.2. La Comprobación de Tipos de las Funciones Polimorfas

#### Reglas a Tener en Cuenta en la Inferencia de Tipos

Las reglas para la comprobación de tipos para funciones polimorfas difiere del de las funciones ordinarias en varios aspectos. Consideremos de nuevo el árbol para la llamada `first(first(q))` decorado con los tipos que deberán ser inferidos:

|                                                                                                        |                                                                                                                                                                                                                  |
|--------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> 1  first : list(ALPHA) -&gt; ALPHA; 2  q : list(list(int)) 3  { 4    first(first(q)) 5  } </pre> | <pre> EXPS(   FUNCTIONCALL[int](     ID[first:F(LIST(INT),INT)],     FUNCTIONCALL[list(int)](       ID[first:F(LIST(LIST(INT)) LIST(INT))],       ID[q:LIST(LIST(INT))]     )   ) # FUNCTIONCALL ) # EXPS </pre> |
|--------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

Del ejercicio de hacer la inferencia para el ejemplo sacamos algunas consecuencias:

1. Distintas ocurrencias de una función polimorfa en el AST no tienen que necesariamente tener el mismo tipo.

En la expresión `first(first(q))` la instancia interna de la función `first` tiene el tipo

`list(list(int)) -> list(int)`

mientras que la externa tiene el tipo

`list(int) -> int`

Se sigue que cada ocurrencia de `first` tiene su propia visión de la variable `ALPHA` que se instancia (se *unifica*) a un valor distinto. Por tanto, en cada ocurrencia de `first`, reemplazaremos la variable de tipo `ALPHA` por una variable *fresca* `ALPHA`. Así habrán variables frescas `ALPHA` para los nodos interno y externo de las llamadas `FUNCTIONCALL(ID[first], ...)`

2. Puesto que ahora las expresiones de tipo contienen variables debemos extender la noción de equivalencia de tipos. Supongamos que una función `f` de tipo `ALPHA -> ALPHA` es aplicada a una expresión `exp` de tipo `list(GAMMA)`. Es necesario *unificar* ambas expresiones. En este caso obtendríamos que `ALPHA = list(GAMMA)` y que `f` es una función del tipo `list(GAMMA) -> list(GAMMA)`
3. Es necesario disponer de un mecanismo para grabar el resultado de la unificación de dos árboles/dags. Si la unificación de dos expresiones de tipo `s` y `s'` resulta en la variable `ALPHA` representando al tipo `t` entonces `ALPHA` deberá seguir representando al tipo `t` conforme avancemos en la comprobación de tipos.

## Programa Árbol para la Inferencia de Tipos

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho$ \
cat -n Trans.trg
1 /***** Polymorphic Type Checking *****/
2 /*
3 eyapp -m Aho::Polymorphism Polymorphism.eyp;
4 treereg -m Aho::Polymorphism Trans.trg
5 */
6
7 {
8 use Aho::Unify qw(:all);
9 }
10
11 functioncall: FUNCTIONCALL($f, $arg)
12 => {
13 my ($name, $line) = @{$f->{attr}};
14
15 my $p = new_var();
16
17 my $q = Parse::Eyapp::Node->hexpand("F", $arg->{t}, $p, \&selfrep);
18
19 if (unify($f->{t}, $q)) {
20 print "Unified\n";
21 print "Now type of ".$f->str." is ".strunifiedtree($q)."\n";
22 }
23 else {
24 die "Type error at $line\n";
25 }
26
27 $FUNCTIONCALL->{t} = representative($p);
28 print "Type of ".$FUNCTIONCALL->str." is ".$FUNCTIONCALL->{t}->str."\n";
```

```

29
30 return 1;
31 }
32
33 tuple: ARGS($x, $y)
34 => {
35 $ARGS->{t} = Parse::Eyapp::Node->hexpand("X", $x->{t}, $y->{t}, \&selfrep);
36 }
37
38 id: ID
39 => {
40 my ($name, $line) = @{$ID->{attr}};
41 our %st;
42 my $ts = $st{$name}->{ts};
43 $ID->{t} = fresh($ts);
44 }

```

- La regla para los identificadores accede al término guardado en la tabla de símbolos que describe el tipo del identificador (atributo `ts`). A partir de él se construye el DAG que representa el tipo llamando a la función `fresh` de `Aho::Unify`. Esta función antes de llamar a `hnew` "refresca" todas las apariciones de variables de tipo. Una variable como `ALPHA` será sustituida por `ALPHA#number`.

```

46 sub fresh {
47 my $ts = shift;
48 my $regexp = shift;
49 $regexp = 'Parse::Eyapp::Node::TYPEVAR::[\w:]+' unless defined($regexp);
50
51 $ts =~ s/($regexp)/1count/g;
52 $count++;
53 my @t = Parse::Eyapp::Node->hnew($ts);
54 representative($_, $_) for @t;
55 return $t[0];
56 }

```

El algoritmo de unificación decora los nodos de los DAGs que están siendo unificados con un atributo (atributo cuyo nombre viene dado por la variable `$set`) que es una referencia al representante de la clase de equivalencia a la que pertenece el nodo. Cada clase tiene un único representante al que denominaremos representante canónico. El representante de un nodo puede obtenerse/modificarse usando la función `representative` proveída por `Aho::Unify`.

```

83 sub representative {
84 my $t = shift;
85
86 if (@_) {
87 $t->{$set} = shift;
88 return $t;
89 }
90 $t = $t->{$set} while defined($t->{set}) && ($t != $t->{$set});
91 die "Representative ($set) not defined!".Dumper($t) unless defined($t->{set});
92 return $t;
93 }

```

Al comienzo del algoritmo cada nodo pertenece a una clase en la que el único nodo es el mismo. Conforme se avanza en el proceso de unificación las clases van aumentando de tamaño.



- La acción asociada con la regla `tuple` correspondiente a expresiones de la forma `x`, `y` en los argumentos es construir el tipo producto cartesiano de los tipos de `x` e `y`.

```

33 tuple: ARGS($x, $y)
34 => {
35 $ARGS->{t} = Parse::Eyapp::Node->hexpand("X", $x->{t}, $y->{t}, \&selfrep);
36 }

```

El método `hexpand` de `Parse::Eyapp::Node` permite construir un DAG compatible con los construidos con `hnew` ya que utiliza la misma cache de memoización. El método admite como último argumento un manejador-decorador de atributos que funciona de manera parecida a como lo hace el correspondiente argumento en `new` y `hnew`. El manejador es llamado por `hexpand` con el nodo que acaba de ser creado como único argumento.

**Ejercicio 13.20.1.** *Podíamos construir el nodo mediante una llamada `Parse::Eyapp::Node->new('X')` y expandiendo el nodo con hijos `$x` y `$y`. Sin embargo, este método no conservaría la propiedad de DAG de los grafos de tipo. ¿Porqué?*

**Ejercicio 13.20.2.** *¿Que ocurriría si el nodo se creara con `Parse::Eyapp::Node->hnew('X')` y luego se expandiera con hijos `$x` y `$y`?*

**Ejercicio 13.20.3.** *¿En que forma sin usar `hexpand` podríamos formar un DAG compatible con `hnew` para el árbol `X($x->{t}->str, $y->{t}->str)`?*

La función `selfrep` - también en `Aho::Unify` - tiene por función iniciar el representante del nodo al propio nodo. De hecho su código es:

```

95 sub selfrep {
96 representative($_[0], $_[0])
97 };

```

- La regla `functioncall` para las llamadas `FUNCTIONCALL($f, $arg)` parte de que el tipo mas general asociado con esta clase de nodo debe tener la forma `F($arg->{t}, $p)` donde `$p` es una nueva variable libre.

```

13 my ($name, $line) = @{$f->{attr}};
14
15 my $p = new_var();
16
17 my $q = Parse::Eyapp::Node->hexpand("F", $arg->{t}, $p, \&selfrep);

```

La función `new_var` (en `Aho::Unify`) proporciona una nueva variable:

```

58 # The name space Parse::Eyapp::Node::TYPEVAR:: is reserved for tree variables
59 # Parse::Eyapp::Node::TYPEVAR::_#number is reserved for generated variables
60 sub new_var {
61 my $p = Parse::Eyapp::Node->hnew("Parse::Eyapp::Node::TYPEVAR::_$count");
62 representative($p, $p);
63 $count++;
64 return $p;
65 }

```

Las variables proveídas tiene la forma `Parse::Eyapp::Node::TYPEVAR::_#number`.

La llamada al método `hexpand` se encarga de crear el DAG asociado con `F($arg->{t},$p)`. El último argumento de la llamada a `hexpand` es el manejador de atributos. Pasamos `&$selfrep` para indicar que se trata de un árbol no unificado.

Acto seguido se pasa a *unificar* el tipo de `$f` y el tipo mas general `F($arg->{t},$p)`.

```
19 if (unify($f->{t}, $q)) {
20 print "Unified\n";
21 print "Now type of ".$f->str." is ".strunifiedtree($q)."\n";
22 }
23 else {
24 die "Type error at line $line\n";
25 }
```

Si la unificación tiene éxito el tipo asociado con la imagen (con el valor retornado) debe ser el tipo unificado de la nueva variable `$p` que es accedido por medio de la función `representative`.

```
27 $FUNCTIONCALL->{t} = representative($p);
```

- Cada vez que la unificación tiene éxito el programa imprime el árbol unificado mediante `strunifiedtree`. La función `strunifiedtree` (en `Aho::Unify`) recorre el árbol de representantes construyendo el término que lo describe:

```
125 sub strunifiedtree {
126 my $self = shift;
127
128 my $rep = representative($self);
129 my @children = $rep->children;
130 my $desc = ref($rep);
131 return $desc unless @children;
132
133 my @d;
134 for (@children) {
135 push @d, strunifiedtree($_);
136 }
137 local $" = ",";
138 $desc .= "(@d)";
139 return $desc;
140 }
```

## Ejemplo

La ejecución del compilador con el programa que hemos usado como ejemplo muestra el proceso de inferencia. Empecemos repasando la estructura del árbol y la representación de las declaraciones:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ \
 usepoly.pl prueba01.ply 2
first type: F(LIST(Parse::Eyapp::Node::TYPEVAR::ALPHA),Parse::Eyapp::Node::TYPEVAR::ALPHA)
q type: LIST(LIST(INT))
first : list(ALPHA) -> ALPHA;
q : list(list(int))
{
 first(first(q))
}
```

```

....
Tree:
EXPS(
 FUNCTIONCALL(
 ID[first],
 FUNCTIONCALL(
 ID[first],
 ID[q]
)
) # FUNCTIONCALL
) # EXPS

```

En primer lugar se calculan los tipos para la llamada interna:

```

Unifying F(LIST(TYPEVAR::ALPHA1),TYPEVAR::ALPHA1) and F(LIST(LIST(INT)),TYPEVAR::_3)
Unifying LIST(TYPEVAR::ALPHA1) and LIST(LIST(INT))
Unifying TYPEVAR::ALPHA1 and LIST(INT)
TYPEVAR::ALPHA1 = LIST(INT)
Unifying LIST(INT) and TYPEVAR::_3
TYPEVAR::_3 = LIST(INT)
Unified
Now type of ID[first] is F(LIST(LIST(INT)),LIST(INT))
Type of FUNCTIONCALL(ID[first],ID[q]) is LIST(INT)

```

Asi se ha inferido que el tipo del argumento de la llamada externa a `first` es `LIST(INT)`. La inferencia para el tipo del nodo externo `FUNCTIONCALL` prosigue con una nueva variable fresca `ALPHA0`:

```

Unifying F(LIST(TYPEVAR::ALPHA0),TYPEVAR::ALPHA0) and F(LIST(INT),TYPEVAR::_4)
Unifying LIST(TYPEVAR::ALPHA0) and LIST(INT)
Unifying TYPEVAR::ALPHA0 and INT
TYPEVAR::ALPHA0 = INT
Unifying INT and TYPEVAR::_4
TYPEVAR::_4 = INT
Unified
Now type of ID[first] is F(LIST(INT),INT)
Type of FUNCTIONCALL(ID[first],FUNCTIONCALL(ID[first],ID[q])) is INT

```

Por último el programa nos muestra el árbol sintáctico abstracto con los tipos inferidos:

```

Tree:
EXPS(
 FUNCTIONCALL[INT](
 ID[first:F(LIST(INT),INT)],
 FUNCTIONCALL[LIST(INT)](
 ID[first:F(LIST(LIST(INT)),LIST(INT))],
 ID[q:LIST(LIST(INT))]
)
) # FUNCTIONCALL
) # EXPS

```

### Un Ejemplo Con Especificación de Tipos Incompleta en C

La expresión `g(g)` en el siguiente programa C muestra la aplicación de una función a ella misma. La declaración de la línea 7 define la imagen de `g` como entera pero el tipo de los argumentos queda sin especificar.

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ \
cat -n macqueen.c
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ cat -n macqueen.c
 1 #include <stdio.h>
 2 #include <stdlib.h>
 3
 4 int n = 5;
 5
 6 int f(int g())
 7 {
 8 int m;
 9
10 m = n;
11 if (m == 0) return 1;
12 else {
13 n = n - 1;
14 return m * g(g);
15 }
16 }
17
18 main() {
19 printf("%d factorial is %d\n", n, f(f));
20 }

```

Al ser compilado con gcc el programa no produce errores ni mensajes de advertencia:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ gcc macqueen.c
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ a.out
0 factorial is 120

```

El error que se aprecia en la salida se corrige cuando cambiamos el orden de los argumentos en la llamada a printf:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ \
sed -ne '/printf/p' macqueen2.c; gcc macqueen2.c ; a.out
printf("%d is the factorial of %d\n", f(f), n);
120 is the factorial of 5

```

¿Cuál es la especificación completa del tipo de g?

Este ejemplo esta basado en un programa Algol dado por Ledgard en 1971 [?]. Aparece como ejercicio 6.31 en el libro de Aho, Sethi y Ullman [?].

El siguiente programa en nuestro pequeño lenguaje polimorfo modela el problema de determinar el tipo de esta función:

```

pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ \
usepoly.pl exercise6_31.ply 2

m : int;
times : int * int -> int;
g : ALPHA
{
 times(m, g(g))
}

```

El árbol y la tabla de símbolos quedan montados después de la primera fase:

Tree:

```

EXPS(
 FUNCTIONCALL(
 ID[times],
 ARGS(
 ID[m],
 FUNCTIONCALL(
 ID[g],
 ID[g]
)
) # ARGS
) # FUNCTIONCALL
) # EXPS
g type: Parse::Eyapp::Node::TYPEVAR::ALPHA
m type: INT
times type: F(X(INT,INT),INT)

```

En primer lugar se procede a la unificación de la llamada interna `FUNCTIONCALL(ID[g],ID[g])`:

```

Unifying TYPEVAR::ALPHA2 and F(TYPEVAR::ALPHA3,TYPEVAR::_4)
TYPEVAR::ALPHA2 = F(TYPEVAR::ALPHA3,TYPEVAR::_4)
Unified
Now type of ID[g] is F(Parse::Eyapp::Node::TYPEVAR::ALPHA3,Parse::Eyapp::Node::TYPEVAR::_4)
Type of FUNCTIONCALL(ID[g],ID[g]) is TYPEVAR::_4

```

Después se pasa a la unificación de la llamada externa `FUNCTIONCALL(ID[times],ARGS(ID[m],FUNCTIONCALL(...`

```

Unifying F(X(INT,INT),INT) and F(X(INT,TYPEVAR::_4),TYPEVAR::_5)
Unifying X(INT,INT) and X(INT,TYPEVAR::_4)
Unifying INT and TYPEVAR::_4
TYPEVAR::_4 = INT
Unifying INT and TYPEVAR::_5
TYPEVAR::_5 = INT
Unified
Now type of ID[times] is F(X(INT,INT),INT)
Type of FUNCTIONCALL(ID[times],ARGS(ID[m],FUNCTIONCALL(ID[g],ID[g]))) is INT

```

Vemos que el tipo de `g` es `F(Parse::Eyapp::Node::TYPEVAR::ALPHA3,Parse::Eyapp::Node::TYPEVAR::_4)` y que `TYPEVAR::_4 = INT`. Así pues el tipo `ALPHA` de `g` es `F(ALPHA,INT)`. Por tanto tenemos que el tipo inferido para `g` es un tipo recursivo:

```
typedef int ALPHA(ALPHA);
```

El programa nos muestra el árbol anotado con los tipos inferidos:

```

Tree:
EXPS(
 FUNCTIONCALL[INT](
 ID[times:F(X(INT,INT),INT)],
 ARGS[X(INT,INT)](
 ID[m:INT],
 FUNCTIONCALL[INT](
 ID[g:F(Parse::Eyapp::Node::TYPEVAR::ALPHA3,INT)],
 ID[g:Parse::Eyapp::Node::TYPEVAR::ALPHA3]
)
) # ARGS
) # FUNCTIONCALL
) # EXPS

```

### Inferencia de una Variable Conocido el Tipo de la Función

La decoración ocurre de abajo-arriba y la inferencia sólo se hace en los nodos de llamada. Sin embargo, Se puede producir la inferencia del tipo de una variable después de que su nodo haya sido visitado, como muestra el ejemplo:

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ \
usepoly.pl reverseinference.ply 2

first : int -> int;
q : ALPHA
{
 first(q)
}
Unifying F(INT,INT) and F(TYPEVAR::ALPHA1,TYPEVAR::_2)
Unifying INT and TYPEVAR::ALPHA1
TYPEVAR::ALPHA1 = INT
Unifying INT and TYPEVAR::_2
TYPEVAR::_2 = INT
Unified
Now type of ID[first] is F(INT,INT)
Type of FUNCTIONCALL(ID[first],ID[q]) is INT
Tree:
EXPS(
 FUNCTIONCALL[INT](
 ID[first:F(INT,INT)],
 ID[q:INT]
)
) # EXPS
first type: F(INT,INT)
q type: Parse::Eyapp::Node::TYPEVAR::ALPHA
```

No obstante, parece razonable que el tipo de una variable que no sea una función debería ser fijo. Dado que las ventajas del polimorfismo se aprecian fundamentalmente en las funciones podemos añadir a nuestro lenguaje reglas que prohíban el polimorfismo de variables que no sean funciones.

#### 13.20.3. El Compilador

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho/script$ cat -n usepoly.p
 1 #!/usr/bin/perl -I../lib -w
 2 use strict;
 3 use Aho::Polymorphism;
 4 use Aho::Unify qw(:all);
 5 use Parse::Eyapp::Base qw(:all);
 6 use Data::Dumper;
 7
 8 # Read Input
 9 my $filename = shift || die "Usage:\n$0 file.c\n";
10 my $debug = shift;
11 $debug = 0 unless defined($debug);
12 my $input = slurp_file($filename, "ply");
13 print $input if $debug;
14
15 # Syntax Analysis and Type Inference
16 my $parser = Aho::Polymorphism->new();
17 my $t = $parser->compile($input);
```

```

18
19 # Show decorated tree
20 $Parse::Eyapp::Node::INDENT = $debug;
21 push_method(qw{FUNCTIONCALL ARGS}, info => sub { strunifiedtree($_[0]->{t}) });
22 push_method(qw{ID}, info => sub { $_[0]->{attr}[0].":".strunifiedtree($_[0]->{t}) });
23 print "Tree: ".$t->str."\n";
24
25 # Print Symbol Table
26 my $symboltable = $t->{symboltable};
27 for (keys(%$symboltable)) {
28 print "$_ type: "
29 . $symboltable->{$_}{ts}
30 ."\n";
31 }

```

### 13.20.4. Un Algoritmo de Unificación

El algoritmo de unificación recibe dos referencias  $\$m$  y  $\$n$  a los árboles que van a ser unificados. Retorna verdadero si la unificación puede tener lugar y falso en caso contrario. La equivalencia se mantiene utilizando un atributo cuyo nombre es  $\$set$ . Todos los nodos en una clase de equivalencia tienen un único representante. Los atributos  $\$set$  de los nodos en una misma clase referencian al representante (posiblemente de forma indirecta, via los campos  $\$set$  de otros nodos). El atributo  $\$set$  del representante canónico apunta a si mismo. Inicialmente cada nodo esta en una clase única.

```

99 sub unify {
100 my ($m, $n) = @_;
101
102 my $s = representative($m);
103 my $t = representative($n);
104
105 return 1 if ($s == $t);
106
107 return 1 if $samebasic->($s, $t);
108
109 print "Unifying ".$s->representative($s)->str." and ".$t->representative($t)->str."\n" if $debug;
110 return 1 if (mergevar($s, $t));
111
112 if (ref($s) eq ref($t)) {
113 $s->{$set} = representative($t);
114 my $i = 0;
115 for ($s->children) {
116 my $tc = $t->child($i++);
117 return 0 unless unify($_, $tc);
118 }
119 return 1;
120 }
121
122 return 0;
123 }

```

Nótese que unificar  $\$m$  y  $\$n$  es unificar sus representantes canónicos  $\$s$  y  $\$t$ . Los representantes serán iguales si  $\$m$  y  $\$n$  ya están en la misma clase de equivalencia.

El algoritmo de unificación usa las siguientes funciones auxiliares:

1. La función **representative** retorna una referencia al representante canónico de la clase:

```

83 sub representative {
84 my $t = shift;
85
86 if (@_) {
87 $t->{$set} = shift;
88 return $t;
89 }
90 $t = $t->{$set} while defined($t->{set}) && ($t != $t->{$set});
91 die "Representative ($set) not defined!".Dumper($t) unless defined($t->{set});
92 return $t;
93 }

```

También permite cambiar el representante.

2. Si `$s` y `$t` representan el mismo tipo básico la función referenciada por `$samebasic` devolverá verdadero. La función `$samebasic` es definida por medio de la función `set` junto con los otros parámetros del algoritmo de unificación.

```

203 Aho::Unify->set(
204 key => 'set',
205 isvar => sub { $_[0] =~ /^Parse::Eyapp::Node::TYPEVAR/ },
206 samebasic => sub {
207 my ($s, $t) = @_;
208
209 return (((ref($s) eq 'INT') || (ref($s) eq 'STRING')) && (ref($s) eq ref($t)));
210 },
211 debug => $debug,
212);

```

3. Si una de los dos árboles `$s` o `$t` es una variable la introducimos en la clase de equivalencia de la otra. Puesto que cada variable fresca es un único nodo (estamos trabajando con un DAG) la actualización de la clase de equivalencia será visible en todas las apariciones de esta variable. Además `representative` recorre la lista de enlaces de unificación retornando el tipo básico o constructor con el que la variable se ha unificado.

La función `mergevar` mezcla las clases de equivalencia cuando uno de los representantes es una variable:

```

67 sub mergevar {
68 my ($s, $t) = @_;
69
70 if (isvar($s)) {
71 $s->{$set} = representative($t);
72 print $s->str." = ".representative($t)->str."\n" if $debug;
73 return 1;
74 }
75 if (isvar($t)) {
76 $t->{$set} = representative($s);
77 print $t->str." = ".representative($s)->str."\n" if $debug;
78 return 1;
79 }
80 return 0;
81 }

```



4. El programador puede proporcionar mediante el método `set` una función `$isvar` que permite añadir variables adicionales:

```
23 sub isvar {
24 my $x = $isvar->(@_);
25 return $x if $x;
26 return 1 if $_[0] =~ /^Parse::Eyapp::Node::TYPEVAR::[\w:]+$\/;
27 }
```

El espacio de nombres `Parse::Eyapp::Node::TYPEVAR` se usa para las variables.

## La Cabecera del Módulo de Unificación

```
pl@nereida:~/doc/casiano/PLBOOK/PLBOOK/code/Aho-Polymorphism/lib/Aho$ cat -n Unify.pm
 1 package Aho::Unify;
 2 use Data::Dumper;
 3 use Parse::Eyapp::Node;
 4 use base qw (Exporter);
 5 our @EXPORT_OK = qw(
 6 unify
 7 new_var
 8 fresh
 9 representative
10 strunifiedtree
11 hnewunifiedtree
12 newunifiedtree
13 selfrep
14);
15 our %EXPORT_TAGS = ('all' => [@EXPORT_OK]);
16
17 my $count = 0;
18 my $set = 'representative';
19 my $isvar = sub { };
20 my $samebasic = sub { };
21 my $debug = 0;
```

## La Función set

```
29 sub set {
30 my $class = shift if @_ %2;
31 $class = __PACKAGE__ unless defined($class);
32
33 my %handler = @_;
34
35 $set = 'representative';
36 $set = $handler{key} if exists($handler{key});
37 $isvar = $handler{isvar} if exists($handler{isvar});
38 $samebasic = $handler{samebasic} if exists($handler{samebasic});
39 $debug = $handler{debug} if exists($handler{debug});
40 $count = 0;
41
42 bless { key => $set, isvar => $isvar, samebasic => $samebasic, count => $count }, $cla
43 }
```

### 13.21. Práctica: Inferencia de Tipos

Extienda el lenguaje con tipos variables presentado en la sección 13.20 con el tipo puntero:

```
q : pointer(string);
```

Suponga añadida la función polimorfa **deref**(**x**) que devuelve lo apuntado por **x**. ¿Cuál es el tipo de **deref**? Extienda las expresiones con constantes y operaciones binarias. Añada código para emitir un mensaje de error en el caso de que después de la fase de inferencia alguna variable - que no sea del tipo función - tenga un tipo polimorfo. Actualice las entradas en la tabla de símbolos para las funciones con un atributo que recoge los diferentes tipos inferidos de los usos de la función.

## Capítulo 14

# Instrucciones Para la Carga de Módulos en la ETSII

Cuando entra a una de las máquinas de los laboratorios de la ETSII encontrará montado el directorio `/soft/perl5lib/` y en él algunas distribuciones de módulos Perl que he instalado.

```
export MANPATH=$MANPATH:/soft/perl5lib/man/
export PERL5LIB=/soft/perl5lib/lib/perl5
export PATH=./home/casiano/bin:$PATH:/soft/perl5lib/bin:
```

Visite esta página de vez en cuando. Es posible que añada algún nuevo camino de búsqueda de librerías y/o ejecutables.

## Capítulo 15

# Usando Subversion

En esta sección indicamos los pasos para utilizar subversión bajo el protocolo SSH. Se asume que ha configurado sub conexión SSH con el servidor de subversion para que admita autenticación automática. Véase El capítulo sobre la SSH en los apuntes de Programación en Paralelo II (<http://nereida.deioc.ull.es/> p para aprender a establecer autenticación automática vía SSH.

**Use Subversion: Creación de un Repositorio** Parece que en banot esta instalado subversion. Para crear un repositorio emita el comando `svnadmin create`:

```
-bash-3.1$ uname -a
Linux banot.etsii.ull.es 2.6.24.2 #3 SMP Fri Feb 15 10:39:28 WET 2008 i686 i686 i386 GNU/Linux
-bash-3.1$ svnadmin create /home/loginname/repository/
-bash-3.1$ ls -l repository/
total 28
drwxr-xr-x 2 loginname apache 4096 feb 28 11:58 conf
drwxr-xr-x 2 loginname apache 4096 feb 28 11:58 dav
drwxr-sr-x 5 loginname apache 4096 feb 28 12:09 db
-r--r--r-- 1 loginname apache 2 feb 28 11:58 format
drwxr-xr-x 2 loginname apache 4096 feb 28 11:58 hooks
drwxr-xr-x 2 loginname apache 4096 feb 28 11:58 locks
-rw-r--r-- 1 loginname apache 229 feb 28 11:58 README.txt
```

Una alternativa a considerar es ubicar el repositorio en un dispositivo de almacenamiento portable (pendriver)

### Añadiendo Proyectos

Ahora esta en condiciones de añadir proyectos al repositorio creado usando `svn import`:

```
[loginname@tonga]~/src/perl/> uname -a
Linux tonga 2.6.24.2 #1 SMP Thu Feb 14 15:37:31 WET 2008 i686 i686 i386 GNU/Linux
[loginname@tonga]~/src/perl/> pwd
/home/loginname/src/perl
[loginname@tonga]~/src/perl/> ls -ld /home/loginname/src/perl/Grammar-0.02
drwxr-xr-x 5 loginname Profesor 4096 feb 28 2008 /home/loginname/src/perl/Grammar-0.02
[loginname@tonga]~/src/perl/> svn import -m 'Grammar Extended Module' \
 Grammar-0.02/ \
 svn+ssh://banot/home/loginname/repository/Grammar
```

```
Añadiendo Grammar-0.02/t
Añadiendo Grammar-0.02/t/Grammar.t
Añadiendo Grammar-0.02/lib
Añadiendo Grammar-0.02/lib/Grammar.pm
Añadiendo Grammar-0.02/MANIFEST
```

```

Añadiendo Grammar-0.02/META.yml
Añadiendo Grammar-0.02/Makefile.PL
Añadiendo Grammar-0.02/scripts
Añadiendo Grammar-0.02/scripts/grammar.pl
Añadiendo Grammar-0.02/scripts/Precedencia.y
Añadiendo Grammar-0.02/scripts/Calc.y
Añadiendo Grammar-0.02/scripts/aSb.y
Añadiendo Grammar-0.02/scripts/g1.y
Añadiendo Grammar-0.02/Changes
Añadiendo Grammar-0.02/README

```

Commit de la revisión 2.

En general, los pasos para crear un nuevo proyecto son:

```

* mkdir /tmp/nombreProyecto
* mkdir /tmp/nombreProyecto/branches
* mkdir /tmp/nombreProyecto/tags
* mkdir /tmp/nombreProyecto/trunk
* svn mkdir file:///var/svn/nombreRepositorio/nombreProyecto -m 'Crear el proyecto nombreProye
* svn import /tmp/nombreProyecto \
 file:///var/svn/nombreRepositorio/nombreProyecto \
 -m "Primera versión del proyecto nombreProyecto"

```

### Obtener una Copia de Trabajo

La copia en Grammar-0.02 ha sido usada para la creación del proyecto, pero no pertenece aún al proyecto. Es necesario descargar la copia del proyecto que existe en el repositorio. Para ello usamos `svn checkout`:

```

[loginname@tonga]~/src/perl/> rm -fR Grammar-0.02
[loginname@tonga]~/src/perl/> svn checkout svn+ssh://banot/home/loginname/repository/Grammar G
A Grammar/t
A Grammar/t/Grammar.t
A Grammar/MANIFEST
A Grammar/META.yml
A Grammar/lib
A Grammar/lib/Grammar.pm
A Grammar/Makefile.PL
A Grammar/scripts
A Grammar/scripts/grammar.pl
A Grammar/scripts/Calc.y
A Grammar/scripts/Precedencia.y
A Grammar/scripts/aSb.y
A Grammar/scripts/g1.y
A Grammar/Changes
A Grammar/README
Revisión obtenida: 2

```

Ahora disponemos de una copia de trabajo del proyecto en nuestra máquina local:

```

[loginname@tonga]~/src/perl/> tree Grammar
Grammar
|-- Changes
|-- MANIFEST
|-- META.yml

```

```

|-- Makefile.PL
|-- README
|-- lib
| '-- Grammar.pm
|-- scripts
| |-- Calc.yp
| |-- Precedencia.yp
| |-- aSb.yp
| |-- g1.yp
| '-- grammar.pl
'-- t
 '-- Grammar.t

```

3 directories, 12 files

```

[loginname@tonga]~/src/perl/>
[loginname@tonga]~/src/perl/> cd Grammar
[loginname@tonga]~/src/perl/Grammar/> ls -la
total 44
drwxr-xr-x 6 loginname Profesor 4096 feb 28 2008 .
drwxr-xr-x 5 loginname Profesor 4096 feb 28 2008 ..
-rw-r--r-- 1 loginname Profesor 150 feb 28 2008 Changes
drwxr-xr-x 3 loginname Profesor 4096 feb 28 2008 lib
-rw-r--r-- 1 loginname Profesor 614 feb 28 2008 Makefile.PL
-rw-r--r-- 1 loginname Profesor 229 feb 28 2008 MANIFEST
-rw-r--r-- 1 loginname Profesor 335 feb 28 2008 META.yml
-rw-r--r-- 1 loginname Profesor 1196 feb 28 2008 README
drwxr-xr-x 3 loginname Profesor 4096 feb 28 2008 scripts
drwxr-xr-x 6 loginname Profesor 4096 feb 28 2008 .svn
drwxr-xr-x 3 loginname Profesor 4096 feb 28 2008 t

```

Observe la presencia de los subdirectorios de control .svn.

## Actualización del Proyecto

Ahora podemos modificar el proyecto y hacer públicos los cambios mediante `svn commit`:

```

loginname@tonga]~/src/perl/Grammar/> svn rm META.yml
D META.yml
[loginname@tonga]~/src/perl/Grammar/> ls -la
total 40
drwxr-xr-x 6 loginname Profesor 4096 feb 28 2008 .
drwxr-xr-x 5 loginname Profesor 4096 feb 28 12:34 ..
-rw-r--r-- 1 loginname Profesor 150 feb 28 12:34 Changes
drwxr-xr-x 3 loginname Profesor 4096 feb 28 12:34 lib
-rw-r--r-- 1 loginname Profesor 614 feb 28 12:34 Makefile.PL
-rw-r--r-- 1 loginname Profesor 229 feb 28 12:34 MANIFEST
-rw-r--r-- 1 loginname Profesor 1196 feb 28 12:34 README
drwxr-xr-x 3 loginname Profesor 4096 feb 28 12:34 scripts
drwxr-xr-x 6 loginname Profesor 4096 feb 28 2008 .svn
drwxr-xr-x 3 loginname Profesor 4096 feb 28 12:34 t
[loginname@tonga]~/src/perl/Grammar/> echo "Modifico README" >> README
[loginname@tonga]~/src/perl/Grammar/> svn commit -m 'Just testing ...'
Eliminando META.yml
Enviando README
Transmitiendo contenido de archivos .

```

Commit de la revisión 3.

Observe que ya no es necesario especificar el lugar en el que se encuentra el repositorio: esa información esta guardada en los subdirectorios de administración de subversion `.svn`

El servicio de subversion parece funcionar desde fuera de la red del centro. Véase la conexión desde una maquina exterior:

```
pp2@nereida:/tmp$ svn checkout svn+ssh://loginname@banot.etsii.ull.es/home/loginname/repositorio
loginname@banot.etsii.ull.es's password:
loginname@banot.etsii.ull.es's password:
A Grammar/t
A Grammar/t/Grammar.t
A Grammar/MANIFEST
A Grammar/lib
A Grammar/lib/Grammar.pm
A Grammar/Makefile.PL
A Grammar/scripts
A Grammar/scripts/grammar.pl
A Grammar/scripts/Calc.yp
A Grammar/scripts/Precedencia.yp
A Grammar/scripts/aSb.yp
A Grammar/scripts/g1.yp
A Grammar/Changes
A Grammar/README
Revisión obtenida: 3
```

## Comandos Básicos

- Añadir y eliminar directorios o ficheros individuales al proyecto

```
svn add directorio_o_fichero
svn remove directorio_o_fichero
```

- Guardar los cambios

```
svn commit -m "Nueva version"
```

- Actualizar el proyecto

```
svn update
```

- Ver el estado de los ficheros

```
svn status -vu
```

- Crear un tag

```
svn copy svn+ssh://banot/home/logname/repository/PL-Tutu/trunk \
 svn+ssh://banot/home/casiano/repository/PL-Tutu/tags/practica-entregada \
-m 'Distribución como fué entregada en la UDV'
```

## Referencias

Consulte <http://svnbook.red-bean.com/> . Vea la página de la ETSII <http://www.etsii.ull.es/svn> . En KDE puede instalar el cliente gráfico KDEsvn.

## **Autenticación Automática**

Para evitar la solicitud de claves cada vez que se comunica con el repositorio establezca autenticación SSH automática. Para ver como hacerlo puede consultar las instrucciones en:

<http://search.cpan.org/~casiano/GRID-Machine/lib/GRID/Machine.pod#INSTALLATION>

Consulte también las páginas del manual Unix de `ssh`, `ssh-key-gen`, `ssh_config`, `scp`, `ssh-agent`, `ssh-add`, `sshd`



# Índice alfabético

- ámbito de la declaración, 632
- ámbito dinámico, 639
- árbol de análisis abstracto, 297, 584
- árbol de análisis sintáctico abstracto, 300
- árbol sintáctico concreto, 279, 450
- árboles, 297, 584
- Benchmark, 436
- LEX, 429
- YYSemval, 408, 548
- bison, 426
- flex, 429
- pos, 91
- yacc, 426
  
- AAA, 297, 584
- Abigail, 260
- abstract syntax tree, 297, 584
- acción de reducción, 391, 538
- acción de transformación árbol, 599
- acción en medio de la regla, 413, 553
- acciones de desplazamiento, 390, 537
- acciones semánticas, 293, 451, 554
- acciones shift, 390, 537
- alfabeto con función de aridad, 297, 584
- algoritmo de construcción del subconjunto, 389, 536
- análisis de ámbito, 631
- antiderivación, 386, 532
- aplicación parcial, 707
- AST, 297, 584
- atributo ámbito, 634
- atributo heredado, 293, 407, 451, 555, 575
- atributo sintetizado, 293, 407, 451, 554, 575
- atributos de los símbolos, 381, 470
- atributos formales, 407, 575
- atributos heredados, 407, 408, 548, 575
- atributos intrínsecos, 407, 575
- atributos sintetizados, 407, 575
- autómata árbol, 317, 591
- autómata finito determinista, 389, 536
- autómata finito no determinista con  $\epsilon$ -transiciones, 388, 534
  
- bloque básico, 333
  
- código auxiliar para las transformaciones, 598
- cabecera, 564
- can, 317
- casa con la sustitución, 316, 591
- casa con un árbol, 316, 591
- casamiento de árboles, 315, 589
- clase, 301
- Class Construction Time, 559
- clausura, 389, 536
- compilador cruzado, 332
- comprobador de tipos, 705
- condición semántica, 599
- conflicto de desplazamiento-reducción, 391, 398, 538, 543
- conflicto reduce-reduce, 392, 398, 538, 543
- conflicto shift-reduce, 391, 398, 538, 543
- conversiones de tipo, 709
- currying, 707
  
- Débilmente Tipado, 706
- DAG, 710
- declaración, 631
- declaración global, 634
- declaración local, 634
- definición dirigida por la sintáxis, 406, 412, 551, 575
- definiciones de familias de transformaciones, 598
- deriva en un paso en el árbol, 298, 585
- DFA, 389, 536
- Directed Acyclic Graph, 710
- Document Type Definition, 153
- documento aquí, 292
- DTD, 153
- duck typing, 709
- dynamic binding, 631
  
- early binding, 631
- Ejercicio
  - Ambigüedad y LL(1), 287
  - Calcular los *FOLLOW*, 285
  - Caracterización de una gramática LL(1), 286
  - Construir los *FIRST*, 284

El orden de las expresiones regulares, 257  
 El or es vago, 258  
 Factores Comunes, 282  
 La opción g, 256  
 Opciones g y c en Expresiones Regulares, 257  
 Recorrido del árbol en un ADPR, 282  
 Regexp para cadenas, 257  
 El nombre de una regla de producción, 571  
 El else casa con el if mas cercano, 357  
 equivalencia de tipos, 653  
 equivalencia de tipos estructural, 708  
 equivalencia de tipos nominal, 708  
 equivalencia por nombres, 708  
 esquema de traducción, 163, 293, 406, 408, 451, 548, 554  
 esquema de traducción árbol, 314, 589  
 Execution Time, 560  
 expresión de tipo, 639, 707  
 expresión regular árbol array, 606  
 expresión regular árbol estrella, 608  
 expresión regular clásica, 600  
 expresión regular lineal, 600  
 expresiones de tipo, 653, 705  
 Expresiones Regulares Arbol, 596  
 expresiones regulares lineales, 601, 720  
 extractores, 345  
 Extreme Programming, 265  
  
 falso bucle for, 566  
 fase de creación del paquete Treeregexp, 602  
 flecha gorda, 600  
 Fuertemente Tipado, 706  
 función de aridad, 297, 584  
 función de transición del autómata, 390, 536  
  
 generador de generadores de código, 588  
 goto, 390, 537  
 grafo de dependencias, 407, 576  
 grafo dirigido acíclico, 710  
 gramática árbol regular, 298, 584  
 gramática atribuída, 407, 576  
 gramática es recursiva por la izquierda, 294  
  
 handle, 387, 533  
 hashed consing, 710  
 here document, 292  
  
 identificación de los nombres, 634  
 idioms, 566  
 Inferencia de Tipos, 737  
 inferencia de tipos, 707  
 inserción automática de anclas, 601  
 inserción automática de la opción x, 601  
 intrínsecamente ambiguos, 354  
 IR, 587  
 isa, 317  
 items núcleo, 393, 540  
  
 L-atribuída, 407, 576  
 LALR, 392, 539  
 late binding, 631  
 lenguaje árbol generado por una gramática, 298, 585  
 lenguaje árbol homogéneo, 297, 584  
 lenguaje generado, 352  
 lenguaje objeto, 706  
 LHS, 571  
 Lisp, 710  
 lista de no terminales, 300, 305  
 LL(1), 286  
 local, 634  
 LR, 386, 532  
  
 máximo factor común, 283  
 método, 301  
 método abstracto, 314  
 manejador de ámbito, 671  
 mango, 387, 533  
 memoization, 711  
 metalenguaje, 706  
 metaprogramas, 706  
 miscreant grammar, 216  
  
 name binding, 631  
 NFA, 388, 534  
 nodo ámbito, 672  
 nodos de uso, 672  
 nombre, 631  
 nombre por defecto, 575  
 normalización del árbol, 315, 589  
 notación dolar, 508, 558, 579  
 notación posicional, 557  
 notación punto, 508, 557, 558, 579  
 notas a pie de árbol, 691  
 nueva opción X, 601  
  
 objeto, 301  
 objeto transformación, 603  
 OCaml, 707  
 ocultar la visibilidad, 635  
 Opción de perl -i, 149  
 Opción de perl -n, 149  
 Opción de perl -p, 149  
 opciones de línea, 149  
 operación de bypass, 502

operaciones de bypass, 502  
orden parcial, 407, 576  
orden topológico, 407, 576

patrón, 315, 589  
patrón árbol, 314, 589  
patrón árbol semántico, 599  
patrón de entrada, 315, 589  
patrón de separación, 349  
patrón lineal, 315, 589  
patrones árbol, 315, 589  
pattern space, 6  
Peephole optimization, 332  
perfilado, 275  
plegado de constantes, 596  
polimorfa, 707  
polimorfismo, 707  
Polimorfismo Ad-hoc, 708  
polimorfismo paramétrico, 708  
polimorfo, 707  
postponed regular subexpression, 126

Práctica

- Ampliación del Lenguaje Simple, 514
- Análisis de Ámbito del Lenguaje Simple C, 652
- Análisis de Tipos en Simple C, 735
- Análisis de Tipos en Simple C con Gramática recursiva por la derecha, 294
- Atribuidas, 735
- Análisis Semántico, 310
- Análisis Sintáctico, 492
- Arbol de Análisis Abstracto, 307
- Autoacciones, 436
- Cálculo de las Direcciones, 321
- Calculadora con Regexp::Grammars, 237
- Casando y Transformando Árboles, 319
- Construcción de los FIRST y los FOLLOW, 285
- Construcción del Arbol para el Lenguaje Simple, 514
- Construcción del AST para el Lenguaje Simple C, 446, 652
- Crear y documentar el Módulo PL::Tutu, 244
- Declaraciones Automáticas, 310
- El Análisis de las Acciones, 436
- Eliminación de la Recursividad por la Izquierda, 296
- Establecimiento de la relación uso-declaración, 700
- Establecimiento de la Relación Uso-Declaración Usando Expresiones Regulares Árbol, 701
- Estructuras y Análisis de Ámbito, 701
- Fases de un Compilador, 252
- Generación Automática de Árboles, 439
- Generación Automática de Analizadores Predictivos, 287
- Generación de Código, 331
- Gramática Simple en Parse::Eyapp, 477
- Inferencia de Tipos, 755
- Números de Línea, Errores, Cadenas y Comentarios, 258
- Nuevos Métodos, 438
- Optimización Peephole, 333
- Plegado de las Constantes, 314
- Pruebas en el Análisis Léxico, 269
- Sobrecarga de Funciones en Simple C, 737
- Traducción de invitation a HTML, 169
- Traductor de Términos a GraphViz, 475
- Un analizador APDR, 287
- Un C simplificado, 421
- Un lenguaje para Componer Invitaciones, 152
- Uso de Yacc y Lex, 430

Primeros, 388, 535  
profiler, 275  
profiling, 275  
programación genérica, 708  
Protocolo YATW de Llamada, 609

recursiva por la izquierda, 294  
reducción-reducción, 392, 398, 538, 543  
reflexividad, 706  
regexp, 600  
regla de producción unaria, 502  
regla por defecto, 25  
reglas de ámbito, 631  
reglas de evaluación de los atributos, 407, 575  
reglas de transformación, 314, 589  
reglas de visibilidad, 635  
reglas semánticas, 407, 575  
rendimiento, 436

Repaso

- Fases de un Compilador, 250
- Las Bases, 243
- Pruebas en el Análisis Léxico, 277

Representación intermedia, 587  
rightmost derivation, 386, 532

Se atribuida, 408, 576  
scope, 631  
script sed, 6  
sección de código, 25  
sección de definiciones, 25  
sección de reglas, 25  
selección de código, 587

siguientes, 388, 535  
sistema de tipos, 705  
SLR, 390, 391, 537, 538  
sobrecarga de identificadores, 639  
sobrecargado, 706  
static binding, 631  
sustitución, 316  
sustitución árbol, 590  
syntactic token, 716  
syntax token, 563

términos, 297, 584  
tabla de acciones, 390, 537  
tabla de gotos, 390, 537  
tabla de saltos, 390, 537  
terminal sintáctico, 563  
test del pato, 709  
tipado dinámico, 705  
tipado estático, 705  
tipado pato, 709  
Tree Construction Time, 559  
Treeregexp, 596  
treeregexp, 598  
trimming, 383, 457

unificar, 744, 747

valores separados por comas, 361  
variables de tipo, 639, 707  
VERSION, 317  
virtual binding, 631

yydebug, 397, 399, 466, 544

zero-width assertions, 94, 117

# Bibliografía

- [1] Dale Dougherty. *Sed and AWK*. O'Reilly & Associates, Inc., Sebastopol, CA, USA, 1991.
- [2] Jerry Peek, Tim O'Reilly, Dale Dougherty, Mike Loukides, Chris Torek, Bruce Barnett, Jonathan Kamens, Gene Spafford, and Simson Garfinkel. *UNIX power tools*. Bantam Books, Inc., New York, NY, USA, 1993.
- [3] Jeffrey E.F. Friedl. *Mastering Regular Expressions*. O'Reilly, USA, 1997. ISBN 1-56592-257-3.
- [4] Casiano Rodriguez Leon. *Perl: Fundamentos*. <http://nereida.deioc.ull.es/~lhp/perlexamples/>, 2001.
- [5] Peter Scott. *Perl Medic: Maintaining Inherited Code*. Addison Wesley, USA, 2004. ISBN 0201795264.
- [6] Ian Langworth and chromatic. *Perl Testing: A Developer's Notebook*. O'Reilly Media, Inc., 2005.
- [7] Damian Conway. *Perl Best Practices*. O'Reilly Media, Inc., 2005.
- [8] Conway D. *Object Oriented Perl*. Manning, Greenwich, USA, 2000.
- [9] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Princiles, Techniques, and Tools*. Addison-Wesley, 1986.
- [10] Todd A. Proebsting. Burg, iburg, wburg, gburg: so many trees to rewrite, so little time (*invited talk*). In *ACM SIGPLAN Workshop on Rule-Based Programming*, pages 53–54, 2002.
- [11] B. Kernighan and D. Ritchie. *The C Programming Language*. Prentice Hall, second edition, 1988.
- [12] Henry F. Ledgard. Ten mini-languages in need of formal definition. *SIGPLAN Not.*, 5(4-5):14–37, 1970.