### Laboratorio

# Procesadores del Lenguaje

2022-2023

### **Presentación**

Profesores Laboratorio.

- Jesús García Herrero
- (jgherrer@inf.uc3m.es)

## **Objetivos**

Plantear una gramática para un lenguaje Verificar si la gramática cumple propiedades para su traducción eficiente

Generar un analizador para un lenguaje dado

Construir los controles semánticos necesarios para verificar y traducir un lenguaje

**Conocer los principios de generación de código** 

## **Desarrollo y Herramientas**

Las Prácticas de Laboratorio:

Se realizarán en grupos de dos

Se programará en

lenguaje Pythonmódulo PLYentorno libre

**Herramientas** 

PLY - lex Scanner
PLY- yacc Parsers

(http://www.dabeaz.com/ply/)

## **Evaluación**

#### Las Prácticas de Laboratorio:

Supondrán un 40% de la nota de la asignatura:

Pla Práctica – Introducción (0.75)

**2**<sup>a</sup> Práctica – Análisis léxico (0.75)

3<sup>a</sup> Práctica – Análisis sintáctico y semántico (2.5)

## INTRODUCCIÓN A LAS HERRAMIENTAS DE COMPILADORES

Procesadores del lenguaje (lex y yacc con PLY)

#### LEX/FLEX

- Analizador léxico (también conocido **Scanner**), escrito originalmente en C, reimplementado en Python en PLY-lex
- o Diseñado para trabajar junto a Yacc, analizador sintáctico para gramáticas LALR.

https://www.cs.princeton.edu/~appel/modern/c/software/flex/flex.html

#### YACC/BISON

- Sirve para generar analizadores sintácticos (o *parsers*) a partir de la especificación de una gramática independiente del contexto.
- YACC/ Bison originalmente desarrollados en lenguaje C, reimplementado en Python en PLY-yacc

https://www.gnu.org/software/bison/manual/bison.html

### Instalación de PLY

- Descargar e instalar PLY:
  - Localizar la última versión: PLY-3.11 http://www.dabeaz.com/ply/
  - Descomprimir la versión a instalar, ej.: ply-3.11.tar.gz
  - Añadir la raíz "python\ply-master" en el path del proyecto
  - Ya pueden importarse los módulos ply.lex, ply.yacc

#### CREACIÓN DEL PROYECTO

- Una vez instalado, podemos crear un nuevo proyecto Python de la forma usual.
  - Por facilidad, en la carpeta "python\example" tenemos algunos ejemplos ilustrativos, en un solo archivo .ply
- En este punto ya podemos empezar a rellenar los archivos de acuerdo a lo especificado en los manuales de lex y yacc

### EJEMPLO PROYECTO (PLY)

```
import sys
sys.path.insert(0, "../..")
# Especificacion de tokens y ERs
# Build the lexer
import ply.lex as lex
lex.lex()
# Parsing rules
# ...
# especificacion de gramática y acciones
# ...
import ply.yacc as yacc
yacc.yacc()
while 1:
    try:
        s = raw_input('calc > ')
    except EOFError:
        break
    if not s:
        continue
    yacc.parse(s)
```

## DEFINIR LENGUAJE LÉXICO: PLANTILLA FLEX

- Realizar las definiciones que sean necesarias para especificar los tokens y expreseiones regulares con la sintaxis de lex
- Reconocer los patrones indicados:
  - Palabras reservadas: lang, es, eng, ...
  - Números
  - Identificadores
  - Comentarios...



Procesadores del lenguaje (PLY-lex)

#### FORMATO DE LEX

- Funcionamiento general: el analizador creado busca en la entrada ocurrencia de los patrones
  - Cuando encuentra un patrón, devuelve el token correspondiente
    - Si tiene función definida, ejecuta sus acciones, y debe devolver el token. Si no, continua la búsqueda
  - Si varios patrones encajan, selecciona el primero declarado
- o Por defecto, si el texto no encaja con ningún patrón se envía a al función de error

## Representación de patrones con Lex

| Literal            | "X"       | La cadena x                 |
|--------------------|-----------|-----------------------------|
| Literal            | \*        | El carácter literal: '*'    |
| Selección          | a ab      | Selección de una            |
|                    |           | alternativa: {a, ab}        |
| Rango              | [ad-gB-   | {a,d,e,f,g,B,C,D,E,G,2,3,4} |
|                    | EG2-4]    |                             |
| Negación del Rango | [^a-z]    | Cualquiera excepto          |
|                    |           | minúsculas                  |
| Agrupar            | (a-z) (0- | Para agrupar patrones [a-   |
|                    | 9)        | z0-9]                       |
| Numeración         | r*        | Ocurrencia de r >=0         |
|                    | r+        | Ocurrencia de r >0          |
|                    | r?        | Cero o una ocurrencia de r  |
|                    | r{2,4}    | De 2 a 4 ocurrencias de r   |
| Cualquier carácter | •         | (. \n)* representa          |
| (no \n)            |           | cualquier fichero           |
| Localización       | ^r        | r al principio de la línea  |
|                    | r\$       | r al final de la línea      |

#### ASPECTOS ADICIONALES PLY-LEX

- Especificacion de tokens con ERs
  - Dar nombre a patrones complejos facilitando la legibilidad
  - Se pueden anidar, otras definiciones se pueden reutilizar

```
tokens = (
    'letter',
    'digit',
    'id',
    'nl',
# tokens DEFINITION
t letter = r'[A-Za-z]'
t digit = r'[0-9]'
t id = r'(' + t letter +
           r'(' + digit + r'|'+ t letter + r')*)'
t_nl = r' n| r| r'
```

#### ASPECTOS ADICIONALES PLY-LEX

• Los tokens pueden especificarse directamente con reglas o con funciones

```
tokens = (
    'NAME', 'NUMBER',
)
# Tokens
t_NAME = r'[a-zA-Z][a-zA-Z0-9_]*'
def t_NUMBER(t)*:
    r'\d+'
    t.value = int(t.value)
    return t
t: clase LexToken
    value
    type
    lineno
    lexpos...
```

• Funciones especiales para **ignorar**, **error**, espacio en blanco

```
t_ignore = " \t"

def t_newline(t):
    r'\n+'
    t.lexer.lineno += t.value.count("\n")

def t_error(t):
    print("Illegal character '%s'" % t.value[0])
    t.lexer.skip(1)
```

• Definición de las **condiciones de arranque** (incl/excl.)

#### ASPECTOS ADICIONALES PLY-LEX

- Condiciones de arranque
  - Para modificar el flujo de análisis: reconocedores más potentes
  - Hay dos tipos de condiciones de arranque
    - Inclusivas
      - o Evaluados los patrones con la condición de arranque y los que no utilizan ninguna
    - Exclusivas
      - Sólo se evalúan los de la condición de arranque
  - Cambiar de estado: función de lex **begin** ()

```
states = (
    ('foo','exclusive'),
    ('bar','inclusive'),
)
t_foo_NUMBER = r'\d+'  # Token 'NUMBER' in state 'foo'
t_bar_ID  = r'[a-zA-Z_][a-zA-Z0-9_]*' # Token 'ID' in state 'bar'

def t_foo_newline(t):
    r'\n'
    t.lexer.lineno += 1
    t.lexer.begin('foo')
def t_bar_newline(t):
    t.lexer.begin('INITIAL')
```

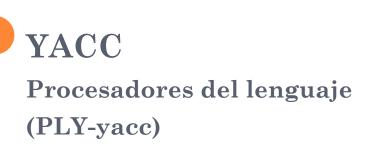
#### Orden de identificación de reglas

- Orden para la identificación de patrones
  - Si para la entrada puedan aplicarse varias reglas:
  - 1. Las reglas literales se ordenan por longitud decreciente
    - Con la entrada abc:

```
t_p1= r'a'
t_p2= r'ab'
t_p3= r'c'
t_p4= r'abc'
```

2. Los patrones definidos por funciones se comprueban en orden de definición

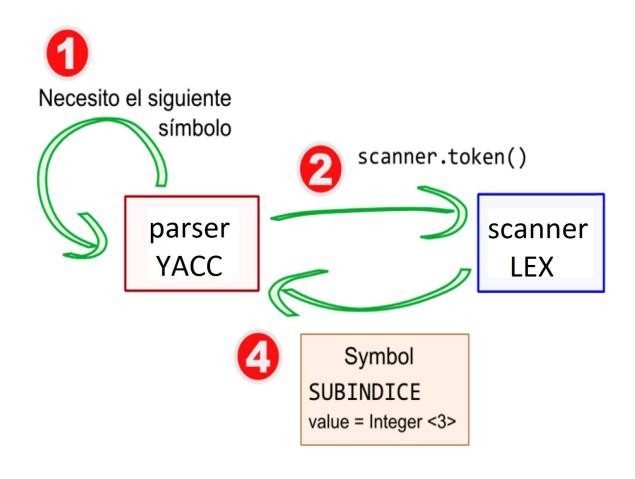
```
def t_REAL(t):
    r'\d+\.\d+'
    t.value = float(t.value)
    return t
def t_ENT(t):
    r'\d+'
    t.value = int(t.value)
    return t
(si cambiamos el orden no funciona)
```



#### Principio de Yacc

- Análisis previo de la gramática para construir el analizador LALR: notifica incidencias
- Funcionamiento general del analizador: pide tokens definidos al analizador léxico
  - Como es LALR(1), pide un token de pre-análisis
  - Cuando realiza desplazamiento pide el siguiente token
  - Cuando realiza reducciones de producciones, ejecuta el código en sus acciones

## COMUNICACIÓN LEX-YACC



## CREACIÓN DEL PROYECTO (YACC)

## Terminales No Terminales

Terminales: {ES, ENG, LANG, SALTO, PALABRA}

no terminales: {S, idioma, texto}

#### Gramática - acciones

S::=LANG idioma texto SALTO
| LANG idioma texto SALTO S

#### DECLARACIONES DE TERMINALES

- Ejemplo definición terminales
  - Lenguaje de expresiones aritméticas con enteros

```
tokens = (
    'NAME', 'NUMBER',
)

literals = ['=', '+', '-', '*', '/', '(', ')']

# Tokens
t_NAME = r'[a-zA-Z_][a-zA-Z0-9_]*'
def t_NUMBER(t):
    r'\d+'
    t.value = int(t.value)
    return t...
```

#### DECLARACIONES DE TERMINALES

• Alternativa definición sin literales

```
tokens = (
    'NAME', 'NUMBER', 'EQUALS', 'PLUS', 'MINUS',
    'BY', 'DIV', 'LBRACE', 'RBRACE',
# Tokens
t NAME = r'[a-zA-Z][a-zA-Z0-9]*'
t EQUALS = r'='
t PLUS = r' + '
t MINUS = r'-'
t_BY = r' \ '*'
t DIV = r' / '
t LBRACE = r' \setminus ('
t RBRACE = r' \setminus )'
```

#### REGLAS EN YACC

Formato BNF simplificado

#### LI: LD acción

- LI: es un símbolo no-terminal del lenguaje
- LD: secuencia de símbolos no-terminales y terminales
- o acción: sentencias en código (puede ser vacío)

### CODIFICACIÓN DE GRAMÁTICA EN PLY

```
# dictionary of names
names = {}
def p statement assign(p);
    'statement : NAME "=" expression'
    names[p[1]] = p[3]
def p statement expr(p):
    'statement : expression'
    print(p[1])
def p expression binop(p):
    '''expression : expression '+' expression
                   | expression '-' expression
                   | expression '*' expression
                   | expression '/' expression'''
    if p[2] == '+':
        p[0] = p[1] + p[3]
    elif p[2] == '-':
        p[0] = p[1] - p[3]
    elif p[2] == '*':
        p[0] = p[1] * p[3]
    elif p[2] == '/':
        p[0] = p[1] / p[3]
```

p[i]: pila semántica. array de objetos del tipo definido para el atributo de cada símbolo de la gramática

#### EJEMPLO DE ANÁLISIS

Gramática

```
expr→expr + term
expr→term
term→term * factor
term→factor
factor→NUMERO
factor→(expr)
```

- o Ejemplo de ejecución
  - 5\*(6+1)+3

```
factor--> NUMERO (5)
term --> factor
factor--> NUMERO (6)
term --> factor
expr --> term
factor--> NUMERO (1)
term --> factor
expr --> expr MAS term
factor--> ( expr )
term --> term POR factor
expr --> term
factor--> NUMERO (3)
term --> factor
```

expr --> expr MAS term

#### RESOLUCIÓN DE AMBIGÜEDAD

- o Por defecto, después de avisar, yacc resuelve
  - D/R: Desplazar prioridad sobre reducir
  - R/R: Reducir por la producción primera
- Preferible resolver los conflictos explícitamente
  - Criterios de prioridad

#### PRECEDENCIA

- Especificación de precedencia
  - Asociatividad izquierda (mismo operador)

```
o ('left', 'op'): x op y op z \rightarrow (x op y) op z
```

Asociatividad derecha (mismo operador)

```
o ('right', 'op'): x op y op z \rightarrow x op (y op z)
```

No asociatividad (mismo operador)

```
o ('nonassoc', 'op'): x op y op z INCORRECTO
```

- Con otros operadores: los declarados en líneas posteriores más precedencia
- Ejemplo de asociación:

```
('left', '+', '-')
('left', '*', '/')
```

• El último declarado es el que tiene más precedencia

# GRAMÁTICA AMBIGÜA: EXPRESIONES CON PRIORIDAD

```
precedence = (
    ('left', '+', '-'),
    ('left', '*', '/'),
    ('right', 'UMINUS'),
def p statement assign(p):
    'statement : NAME "=" expression'
    names[p[1]] = p[3]
def p statement expr(p):
    'statement : expression'
    print(p[1])
def p expression binop(p):
    '''expression : expression '+' expression
                   | expression '-' expression
                   | expression '*' expression
                   | expression '/' expression'''
```

# GRAMÁTICA AMBIGÜA: EXPRESIONES CON PRIORIDAD

```
precedence = (
                                   calc > 2*3+1+2+8
    #('left', '+', '-'),
                                    28
    #('left', '*', '/'),
                                   calc >
    ('right', 'UMINUS'),
def p statement assign(p):
    'statement : NAME "=" expression'
    names[p[1]] = p[3]
def p statement expr(p):
    'statement : expression'
    print(p[1])
def p expression binop(p):
    '''expression : expression '+' expression
                   | expression '-' expression
                   | expression '*' expression
                   | expression '/' expression'''
```

# GRAMÁTICA AMBIGÜA: EXPRESIONES CON PRIORIDAD

```
precedence = (
                                   calc > 2*3+1+2+8
    ('left', '+', '-'),
    ('left', '*', '/'),
                                   calc >
    ('right', 'UMINUS'),
def p statement assign(p):
    'statement : NAME "=" expression'
    names[p[1]] = p[3]
def p statement expr(p):
    'statement : expression'
    print(p[1])
def p expression binop(p):
    '''expression : expression '+' expression
                   | expression '-' expression
                   | expression '*' expression
                   | expression '/' expression'''
```